

POLIREMATICHE E COLLOCAZIONI DELL'ITALIANO

UNO STUDIO LINGUISTICO E COMPUTAZIONALE

Luigi Squillante

luna di miele
capelli castani
messaggio
fare i natali
prestare attenz

Impressum

Dieses Werk ist mit der Creative-Commons-Nutzungslizenz «Namensnennung – Nicht kommerziell – Keine Bearbeitung 3.0 Deutschland» versehen. Weitere Informationen finden sind unter: <http://creativecommons.org/licenses/by-nc-nd/3.0/de>

Universitätsverlag Hildesheim
Universitätsplatz 1
31141 Hildesheim

www.uni-hildesheim.de/bibliothek/universitaetsverlag-open-access

Erstausgabe Hildesheim 2016
Redaktion & Satz: Luigi Squillante
Titelbildgestaltung: Mario Müller

ISBN-A 10.978.3934105/720



SAPIENZA
UNIVERSITÀ DI ROMA

FACOLTÀ DI LETTERE E FILOSOFIA
DIPARTIMENTO DI SCIENZE DOCUMENTARIE, LINGUISTICO-FILOLOGICHE E
GEOGRAFICHE

DOTTORATO DI RICERCA IN
FILOLOGIA, LINGUISTICA E LETTERATURA
XXVII CICLO

POLIREMATICHE E COLLOCAZIONI DELL'ITALIANO.
UNO STUDIO LINGUISTICO E COMPUTAZIONALE

TESI SVOLTA IN COTUTELA CON LA
UNIVERSITÄT HILDESHEIM

RELATORI
Prof. Isabella Chiari
Prof. Ulrich Heid

CANDIDATO
Dott. Luigi Squillante
MATRICOLA 1098620

Anno Accademico 2013-2014

Indice

Sigle e abbreviazioni	v
Introduzione	vii
1 Parole di più parole	1
1.1 La questione della parola e il livello del lessico	1
1.2 Le anomalie delle espressioni di più parole	4
1.3 La normalità delle espressioni di più parole	8
1.4 Polirematiche e collocazioni	10
1.5 Approcci classici ai fenomeni idiomatico-collocativi	12
1.5.1 La visione strutturalista	13
1.5.2 La visione generativista	17
1.5.3 Critiche al modello generativo e interpretazioni pragmatiche .	19
1.5.4 Livello del lessico e collocazioni	21
2 L'approccio della linguistica computazionale	25
2.1 Calcolabilità, modelli di lingua, corpora	25
2.2 Principali interessi applicativi	30
2.3 Collocazioni empiriche	32
2.4 Il trattamento informatico e statistico delle parole	35
2.4.1 Definizioni e concetti generali	35
2.4.2 Frequenze di cooccorrenza	39
2.4.3 Misure d'associazione	41
2.4.4 Valutazione dei risultati	46
2.4.5 Limiti delle misure d'associazione	47
3 Uno strumento per le analisi variazionali dei fenomeni multiparola	49
3.1 Approcci computazionali e categorizzazione	49
3.2 Motivazioni per uno studio della modificabilità empirica	53
3.3 Metodologia di studio delle variabilità	54
3.3.1 Corpus, pattern e lista di espressioni	55
3.3.2 Variazioni sintagmatiche	56
3.3.3 Variazioni paradigmatiche	60
3.3.4 Variazioni flessive	69
3.4 Una scelta metodologica sulla valutazione della categorizzazione . . .	70

4	Analisi sul linguaggio generale dell'italiano	73
4.1	Il corpus PAISÀ	73
4.2	Analisi sui pattern nominali	75
4.2.1	Analisi sul pattern NA	82
4.2.2	Analisi sul pattern AN	87
4.2.3	Analisi sul pattern NPN	92
4.2.4	Analisi sul pattern NPdN	97
4.2.5	Analisi sul pattern NPV _{inf}	101
4.2.6	Analisi sul pattern NN	105
4.2.7	Analisi sul pattern NCN	108
4.2.8	Analisi sul pattern VCV	113
4.3	Analisi sul pattern verbale VDN	119
4.4	Conclusioni	123
5	Studio sul pattern NA: il caso della fisica	125
5.1	Terminologia del linguaggio fisico	125
5.2	Il corpus	128
5.2.1	Costituzione del corpus	128
5.2.2	Trattamento	131
5.3	Caratteristiche delle polirematiche fisiche NA	135
5.4	Indice di prototipicità	137
5.5	Analisi e risultati	138
5.6	Conclusioni	141
	Conclusioni e lavori futuri	143
A	Part of speech Tagset - PAISÀ	145
B	Dati sul pattern NA	149
C	Dati sul pattern AN	163
D	Dati sul pattern NPN	177
E	Dati sul pattern NPdN	191
F	Dati sul pattern NPV_{inf}	205
G	Dati sul pattern NN	219
H	Dati sul pattern NCN	233
I	Dati sul pattern VCV	247
J	Dati sul pattern VDN	261

K Part of speech Tagset - Fisica	275
Bibliografia	277

Sigle e abbreviazioni

A	=	aggettivo
Avv	=	avverbio
C	=	congiunzione
CLI	=	clitico
D	=	articolo determinativo
D _{ind}	=	articolo indeterminativo
N	=	sostantivo
N _{pr}	=	nome proprio
NUM	=	aggettivo numerale
P	=	preposizione semplice
Pd	=	preposizione articolata
Pron	=	pronome
V	=	verbo
V _{aux}	=	verbo ausiliare
V _{inf}	=	verbo infinito
V _{pp}	=	verbo al participio passato

Introduzione

Nella linguistica contemporanea, la definizione di un certo insieme di entità a cui ci si riferisce generalmente attraverso una vasta e varia terminologia (*polirematiche, collocazioni, lessemi complessi, composti sintagmatici, espressioni multiparola*, ecc.) rimane controversa. Risulta intuitivamente chiaro che alcune parole, quando appaiono congiunte, esibiscono un “legame peculiare” che può, di volta in volta, definirle ad esempio come unità semantiche (*pollice verde*), sintattiche (*per mezzo di*), come combinazioni che mostrano preferenzialità lessicali (*disputare un campionato*). Ciononostante non appare chiaro il confine lungo il *continuum* lessicosintattico che riesca a separare questo tipo di entità da ciò a cui ci si riferisce generalmente come ‘espressioni libere’; né esiste una netta categorizzazione che possa distinguere le diverse tipologie di espressioni in base a proprietà o comportamenti chiaramente identificabili come pertinenti a sottogruppi definiti.

Il presente lavoro mira a fornire una nuova prospettiva nell’ambito della categorizzazione di tale tipo di espressioni mediante un loro trattamento computazionale attraverso l’uso di *corpora*, tendente a una sistematizzazione che scaturisca dal diverso comportamento delle entità rispetto a parametri linguistici e statistici, in una prospettiva *bottom-up*. La lingua oggetto dello studio è l’italiano, carente di lavori che indaghino sistematicamente a livello quantitativo il comportamento empirico delle espressioni in esame.

La tesi si compone di cinque capitoli.

Il **Capitolo 1** è dedicato alla definizione e all’inquadramento dell’oggetto di studio del lavoro, con una focalizzazione sul rapporto tra le espressioni studiate e il concetto di parola, per poi analizzare le caratteristiche di anomalia e di normalità che tali entità esibiscono rispetto all’intuitivo concetto di comportamento semantico e morfosintattico ‘standard’. L’ultima parte del capitolo è infine dedicata ad una breve panoramica sugli approcci che gli studi linguistici del XX secolo hanno avuto nei confronti delle espressioni multiparola.

Il **Capitolo 2** si focalizza, nello specifico, sull’approccio della linguistica computazionale alle espressioni in esame. Il discorso spazia quindi da un’introduzione su motivazioni e utilità degli studi computazionali, ai principali interessi applicativi del trattamento informatico delle espressioni multiparola, alla definizione dei concetti e delle metodologie più strettamente legati all’interazione tra l’apparato matematico-statistico e la lingua.

Il **Capitolo 3** è dedicato alla descrizione di uno strumento computazionale costruito *ad hoc* per il presente studio e che costituisce il fulcro dell’apparato metodologico della tesi. Grazie a tale strumento è possibile, infatti, eseguire automaticamente una serie di test empirici sui comportamenti *variazionali* di vasti insiemi di espres-

sioni, una volta che si abbia un corpus a disposizione. Le variazioni studiate sono di tipo sintagmatico, paradigmatico e morfologico. Dall'analisi dei dati scaturiti dai test è possibile, quindi, tentare di proporre una categorizzazione delle espressioni in base alle caratteristiche che le espressioni esibiscono una volta che siano state raggruppate in insiemi omogenei rispetto a determinati parametri.

Il **Capitolo 4** descrive, quindi, le analisi e i risultati effettuati grazie allo strumento computazionale su diversi insiemi di espressioni relative a specifici pattern nominali e verbali. I dati sono prodotti grazie all'utilizzo di un vasto corpus (PAISÀ, Lyding *et al.* 2014), rappresentativo dell'italiano generale, che presenta un alto livello di annotazione. Al termine dell'analisi su ogni pattern considerato, viene proposta una categorizzazione delle espressioni secondo diverse tipologie che possono variare a seconda della sequenza grammaticale considerata.

Il **Capitolo 5**, infine, è incentrato su un caso di studio circoscritto al linguaggio tecnico-specialistico della fisica. Sono qui analizzate le proprietà delle espressioni nome-aggettivo costituenti la terminologia del settore, grazie ai test dello strumento computazionale su un corpus raccolto appositamente per lo studio e descritto in tutte le fasi di processamento e annotazione. La parte finale del capitolo è invece dedicata all'esplorazione della possibilità di utilizzare le informazioni variazionali sulle espressioni ai fini dell'individuazione e dell'estrazione automatica da corpora della terminologia di settore.

1

Parole di più parole

1.1 La questione della parola e il livello del lessico

Nella tradizione degli studi linguistici la definizione di parola è da lungo tempo oggetto di riflessione e dibattito a causa della problematicità di individuazione delle caratteristiche necessarie e sufficienti a identificare il fenomeno¹. Il concetto di cosa sia una parola è intuitivamente presente nella coscienza dei parlanti, tuttavia, com'è noto, una definizione "scientifica" e universale non sembra possibile sia a causa dei diversi livelli di analisi cui la parola può essere sottoposta, sia per la varietà dei tipi di lingue esistenti. Saussure sottolinea storicamente la questione, affermando che «le mot, malgré la difficulté qu'on a à le définir, est une unité qui s'impose à l'esprit, quelque chose de centrale dans le mécanisme de la langue» (Saussure, 1922, p. 154). In termini assiomatici, si può dire che la parola sia un "primitivo" della teoria linguistica, «una nozione sulla quale c'è accordo intuitivo, ma per la quale non è possibile approfondire l'analisi» (Scalise & Bisetto, 2008, p. 62). L'idea di unità che sta dietro il concetto di parola è infatti relativa, in quanto lo stesso materiale linguistico può, allo stesso tempo, costituire un *unicum* o un insieme di parti a seconda della prospettiva che si consideri. Grandi (2006, p. 31) ammette che «i numerosi tentativi di produrre una definizione univoca, condivisa e fondata su basi teoricamente solide non hanno sortito effetti di rilievo»². Come ipotizza Ramat (1990), il concetto di parola sarebbe meglio rappresentabile in maniera prototipica, con una struttura centrale che comprenda le entità che manifestano pienamente caratteristiche di coesione interna, espressione di un significato, mobilità e isolabilità, e una periferia che includa gradualmente, in maniera proporzionale alla distanza dal centro, tutte le entità che progressivamente risultano non manifestare *in toto*

¹Per approfondire nello specifico la definizione di parola si vedano i lavori di sintesi di Di Sciullo & Williams (1987), Lepschy (1989), Simone (1990), Ramat (1990, 2005).

²Voghera, tuttavia, puntualizza come sulla questione della parola, da lei definita "annosa", si sia giunti ad una soluzione «ampiamente condivisa», secondo cui esistono dei criteri efficaci in grado di identificare la parola, sebbene abbiano «pregnanza diversa»: *non interrompibilità* (non è possibile inserire materiale sintagmatico all'interno di una parola), *non mobilità dei costituenti* (non è possibile modificare l'ordine dei morfemi), *isolabilità* (possibilità di costituire autonomamente un enunciato), *pausa potenziale* (è sempre possibile inserire prima e dopo una pausa). Di questi, solo i primi due risulterebbero effettivamente validi (Voghera, 1994, p. 190).

le proprietà appena menzionate. Come accennato, il problema ha una delicatezza intrinseca legata alla prospettiva interlinguistica. Ciò che viene considerato parola in una lingua, può non corrispondere ad una parola in un'altra ma, soprattutto, lingue isolanti e agglutinantanti risultano ad esempio agli estremi nel problema di quanto le unità percepite dai parlanti siano effettivamente parole. Sul piano intralinguistico, invece, benché la nozione di parola sia centrale ad ogni livello di analisi, essa non risulta invariante tra i vari livelli, e in ognuno si necessita di una definizione specifica che individui le proprietà più strettamente pertinenti alla prospettiva in esame³. Spesso, inoltre, è la scrittura a corroborare nei parlanti l'idea di una classe astratta di unità distinte, nonostante la sequenza del parlato non mostri interruzioni. Come nota Ramat (2005, p. 107): «quasi tutte le tradizioni scritte [...] conoscono la divisione in parole della sequenza sonora». In particolare, nella nostra tradizione, la parola è spesso stata fatta coincidere con l'unità compresa tra due spazi bianchi in un testo scritto⁴, parlando in questo caso di *parola grafica*.

Sul piano semantico e lessicale esistono espressioni che, nonostante siano formate da più di una parola grafica, possono apparire come un'unità o mostrano legami specifici che non rendono le parole costituenti completamente libere. L'idea più intuitiva di tali manifestazioni linguistiche è rappresentata da espressioni quali *luna di miele*, *colpo di stato*, *pollice verde*, in cui l'unione di più parole grafiche è necessaria alla creazione di un'unità più grande che, nel caso dei tre esempi citati, è identificabile sul piano semantico. Ma allargando appena il campo di vista nella ricerca di espressioni che ad un qualsiasi livello di analisi linguistica appaiano o si comportino come unitarie, nonostante siano formate da più parole, è possibile imbattersi in locuzioni del tipo *fantanto che*, *alla bell'e meglio* o anche *delitto efferato*, che si discostano, in qualche modo, dalle caratteristiche degli esempi citati sopra. La prima differenza percepibile è che *fantanto che* è un'unità sintattica più che semantica, riconducibile alla categoria di congiunzione; *alla bell'e meglio* mostra invece una costruzione sintattica anomala e viene a configurarsi come una locuzione avverbiale di modo. Il discorso è diverso per l'ultimo esempio, dove il legame che intercorre tra i due costituenti sembra essere più debole, in quanto l'occorrenza congiunta delle parole non è utile alla creazione di un nuovo significato né di un componente sintattico, tuttavia *efferato* raramente è usato in contesti in cui non sia presente anche *delitto*

³Sul piano fonologico, ad esempio, si parla di parola (fonologica) intendendo l'insieme delle sillabe che si raccolgono intorno all'accento, come ad esempio *per favore*, in cui le due entità che si considerano generalmente distinte, la preposizione e il sostantivo, sono parte di un *unicum*. Ad un livello più astratto, invece, si considera parola (lessema) l'insieme di tutte le forme flesse di cui una è convenzionalmente scelta a rappresentarle (come nel caso di *correvamo*, *corro*, *correndo*, tutte varianti del lessema *correre*).

⁴È d'obbligo sottolineare che per *nostra tradizione* si intende la tradizione scritta occidentale moderna, in quanto esistono esempi di lingue classiche scritte in cui non esisteva divisione tra le parole (greco, latino), in cui la divisione consisteva in un punto (latino), e lingue attuali con scrittura senza spazi (cinese, giapponese). È utile ricordare come esistano anche lingue senza sistemi di fissazione grafematica a cui è inapplicabile un discorso di parola identificata in base alle convenzioni di scrittura.

e nell'uso è diventato un suo attributo preferenziale.

In generale l'insieme di queste espressioni è molto vasto e dai contorni sfumati. Essi si collocano infatti su un *continuum* graduale, in una zona intermedia tra il lessico e la sintassi, risultando spesso devianti da entrambi. La problematicità stessa di definizione e individuazione di tali fenomeni ha favorito la nascita di un gran numero di termini, come *lessema complesso*, *polirematica*, *espressione idiomatica*, *frase fissa*, *collocazione*, *locuzione*, *frasema*, *fraseologia*, *modo di dire*, ecc. le cui differenze sono anch'esse poco chiare o sfumate; tutti, comunque, indicano espressioni che appaiano in qualche modo cristallizzate o irrigidite nell'uso.

Nonostante la tradizione linguistica abbia sempre avuto ben presente l'esistenza del fenomeno, esso è stato spesso trascurato dagli studiosi. Nella prefazione al testo di Casadei (1996, p. III) Simone, in merito ad esempio alle espressioni idiomatiche, afferma:

i pochi [linguisti] che se ne sono occupati [...] hanno liquidato le espressioni idiomatiche come entità monolitiche sintatticamente e soprattutto non articolabili semanticamente: una pietra d'inciampo, quindi, per la teoria linguistica, che è una dottrina primordialmente analitica, e mal sopporta di incontrare entità che non siano sottoponibili a scomposizione.

Bosque (2004a), in merito alle associazioni preferenziali di parole note come *collocazioni* parla di un "terreno de nadie" negli studi linguistici, in quanto in generale né al grammatico, né al fraseologo, né al lessicografo spettava, all'interno delle proprie mansioni, lo studio di espressioni che non mostravano restrizioni categoriali generalizzabili, né carattere locuzionale e che non avevano bisogno di definizioni a sé stanti, in quanto non costituivano espressioni fisse⁵.

In generale le ragioni per cui l'ambito di espressioni idiomatiche, lessemi complessi, collocazioni e delle "espressioni di più parole" è rimasto a lungo un terreno confuso e inesplorato sono molteplici, varie e difficilmente individuabili nello specifico. Da un lato, il fatto che espressioni rientranti in questo insieme rappresentassero una scomoda anomalia per molte teorie (in virtù della loro già citata devianza dalla sintassi e dalla semantica "standard" della lingua) ha favorito la diffusa e fallace percezione di una loro marginalità all'interno del sistema-lingua. Dall'altro è stato necessario tempo e sviluppo tecnologico per avere a disposizione grandi basi dati testuali (i cosiddetti *corpora*) per attestare l'importanza del livello lessicale e delle preferenze di cooccorrenza delle parole nelle analisi linguistiche.

In ogni caso, come si vedrà nel seguito, vari sono stati gli approcci a questo tipo di fenomeni provenienti da diversi filoni delle teorie sul linguaggio durante tutto il XX secolo. In particolare, molti degli studi sui lessemi complessi hanno visto

⁵Bosque cita il caso dell'avverbio *profundamente* che ben si accosta a verbi come *dormir*, *influir*, ma non a *caber* o *preguntar*. In questo caso non sembra esistere, se non sul piano dell'uso lessicale, una serie di caratteristiche utili a discriminare la categoria di verbi accettabili. Al contempo *dormir profundamente* non sembra possedere lo status di locuzione che gli garantirebbe una definizione dedicata in un qualche dizionario lessicale.

un'impennata negli ultimi decenni grazie alla spinta fornita dalle nuove applicazioni informatiche di trattamento dei dati. Ciononostante, gli studi non hanno prodotto, ad oggi, un'univoca ed esaustiva sistematizzazione del fenomeno.

1.2 Le anomalie delle espressioni di più parole

L'idiomaticità delle espressioni che mostrano un legame che travalica i confini della parola grafica scaturisce spesso dal loro diverso comportamento linguistico rispetto alle consuetudini morfosintattiche e semantiche caratteristiche delle parole semplici di una lingua. Tale devianza può manifestarsi sia nel comportamento dell'espressione come unità (in opposizione paradigmatica alle altre parole) che in quello dei suoi costituenti. Le anomalie sembrano sottolineare la natura ambivalente e intermedia delle espressioni di più parole, che assumono, a seconda dei casi, tratti che normalmente sono caratteristici della parola semplice, del sintagma o della frase. Nonostante le molte variazioni sul tema avute nel corso del tempo, tradizionalmente è possibile ricondurre le anomalie idiomatiche a quattro categorie principali, vale a dire la *agrammaticalità*, la *non sostituibilità*, la *non modificabilità* e la *non composizionalità*, riguardanti rispettivamente grammatica, lessico, morfosintassi e semantica. Va tenuto a mente che tali proprietà non sono condizioni necessarie e sufficienti, ma rappresentano tendenze e caratteristiche possibili di anomalia. Esse possono quindi cooccorrere congiuntamente o singolarmente e presentarsi in vari gradi.

La *agrammaticalità* si esplica in espressioni cristallizzate che presentano una struttura sintattica mal formata nonostante siano perfettamente comprensibili. Esempi italiani sono il già citato *alla bell'e meglio* o *essere in forse*. Fenomeni di questo tipo non sembrano essere presenti in gran numero nella lingua e si attribuisce alla loro costruzione anomala la peculiare non calcolabilità del significato globale (De Mauro & Voghera, 1996).

Per quanto riguarda la *non sostituibilità*, essa si manifesta nella caratteristica *fissità* lessicale (più o meno marcata) che non permette di sostituire alcuni componenti con altri che condividano le stesse caratteristiche categoriali⁶, a meno di non perdere il senso idiomatologico dell'espressione. Un esempio è l'impossibilità di sostituire *colonna sonora* con **pilastro sonoro*. Sembra, quindi, che il legame di parola travalichi il confine dei singoli costituenti, riducendone l'autonomia, ed essi vengano a costituirsi quasi come morfemi, i quali, se sostituiti, modificano o invalidano il significato originale. È possibile incontrare anche casi in cui, indipendentemente dalla presenza di un significato idiomatologico, la sostituzione del costituente, benché possibile, restituirebbe un'espressione innaturale per un parlante nativo: si pensi al caso di *capelli castani* rispetto al comprensibile ma inusuale *capelli marroni*. La *non sostituibilità* conduce quindi a una critica dell'analisi componenziale a meno di non considerare le

⁶Si intendono, qui, parole che abbiano una vicinanza semantica tale da permetterne l'interscambiabilità in molti contesti e che quindi presentino un alto grado di sinonimia.

espressioni di più parole come unità monolitiche inanalizzabili (cosa discutibile per *capelli castani*). Specie nell'ultimo esempio, infatti, non sembrano esserci pertinenti proprietà di restrizione che possano efficacemente discernere tra *castano* e *marrone* in base a tratti specifici del sostantivo determinato; inoltre, come già sottolineava Halliday (1966), emerge chiara la necessità di considerare il lessico come un livello dell'analisi linguistica al fine di tenere conto di relazioni sintagmatiche di questo tipo.

La non modificabilità morfosintattica risulta invece una macro-proprietà che ingloba diversi comportamenti. Alcune espressioni presentano, infatti, un *ordine non modificabile*, non tollerando l'inversione dei loro costituenti (come in **risposta e botta* invece del corretto *botta e risposta*) mentre per altre espressioni vale la *non interrompibilità*, secondo cui non è possibile inserire materiale linguistico tra i costituenti (si pensi a **luna meravigliosa di miele*). Anche qui l'analogia con la parola è pregnante: la non modificabilità e la non interrompibilità delle espressioni sono analoghe alle restrizioni morfemiche all'interno della parola semplice. Infine la *flessione bloccata* è un'altra delle proprietà di non modificabilità possibili per le espressioni idiomatiche, secondo cui le parti normalmente variabili del discorso si cristallizzano perdendo la propria flessione⁷ (è possibile *alti e bassi* ma non **alto e basso*).

Anche per quanto riguarda le trasformazioni sintattiche le espressioni di più parole mostrano anomalie. Per alcune non è ammessa la pronominalizzazione (ad es. per *colpo di stato* → **è quello di stato il colpo a cui mi riferivo*), la topicalizzazione (**è di stato il colpo?*) e per alcune espressioni idiomatiche verbali sono impossibili la passivizzazione (ad es. *vuotare il sacco* → **il sacco è stato vuotato*), l'interrogazione (*cosa ha vuotato Marco?* - *Il sacco*) e la nominalizzazione dell'azione (*tirare le cuoia* → **il tiro delle cuoia*).

Fraser (1970) ha proposto una categorizzazione delle espressioni idiomatiche sulla base della loro modificabilità sintattica, la cosiddetta *frozenness hierarchy*, composta da 7 livelli che esprimono possibili operazioni sintattiche, ordinati in modo che ogni livello erediti tutte le caratteristiche dei livelli inferiori:

Livello 6: espressioni idiomatiche completamente libere;

Livello 5: ricostituzione in un'altra organizzazione strutturale dei costituenti (es. nominalizzazione dell'azione);

⁷Sulla flessione bloccata è interessante l'osservazione di De Mauro & Voghera (1996) riguardo lessemi complessi nominali e verbali. In italiano sia i nomi che i verbi presentano un paradigma di flessione ma, qualora un verbo figuri in un lessema complesso, esso ha grande probabilità di perdere il suo statuto categoriale, poiché la flessione bloccata non è compatibile con le caratteristiche di definizione che sua la categoria impone (persona, tempo, aspetto, modo). Se si vuole, la flessione bloccata impedisce la saturazione delle valenze (a meno degli impersonali, come *può darsi*). Questo fa sì che i verbi siano particolarmente restii a perdere la propria flessione, ma, nei casi in cui ciò accade, essi subiscono un processo di transcategorizzazione per cui *fai da te* o *cessate il fuoco* diventano sostantivi. I nomi, invece, rientrando nell'insieme grammaticalmente presente degli invariabili, possono conservare il proprio statuto categoriale.

- Livello 4: estrazione di un costituente in posizione esterna all'espressione idiomatrica (es. il passivo);
- Livello 3: permutazione di due costituenti adiacenti;
- Livello 2: inserzione di un costituente;
- Livello 1: aggiunta di costituenti non idiomatichi;
- Livello 0: espressioni totalmente immodificabili.

Studi successivi, come quello di Cutler (1982), hanno evidenziato che le espressioni tendono a collocarsi in livelli di maggiore rigidità della gerarchia quanto più lontana nel tempo è la loro apparizione nella lingua. Risultati di questo tipo sembrano quindi indicare che l'idiomaticità favorisca un processo di delessicalizzazione che in generale può comportare la perdita dei significati dei costituenti a favore del significato globale dell'espressione.

Si giunge, quindi, alla principale caratteristica che sancisce l'irregolarità di molte espressioni di più parole sul piano semantico e che in generale è sempre apparsa come il principale indicatore di idiomaticità, vale a dire la violazione del *principio di composizionalità*, secondo cui il significato di un'espressione è unicamente funzione dei significati delle sue parti e delle regole sintattiche con cui esse si combinano.

La non composizionalità può presentarsi in diversi gradi. Esempi come *tirare le cuoia* o *pan di Spagna* sono casi in cui il significato globale non è ricostruibile a partire da quello dei costituenti, non è cioè somma diretta dei singoli significati che compongono l'espressione. Esistono però anche situazioni in cui nel significato globale sono presenti i significati dei suoi costituenti con in aggiunta un non precisato "sovrappiù semantico" che viene a crearsi solo quando tutti i costituenti occorrono insieme. È questo il caso, ad esempio, di *camera oscura* o *cerchio in lega*: non tutte le camere poco illuminate o buie sono utilizzate dai fotografi, eppure l'occorrenza congiunta di *camera* e *oscura* viene necessariamente a indicare, nell'uso, il tratto semantico aggiuntivo "adibita allo sviluppo delle pellicole fotografiche". Analogamente possono esistere dei cerchi forgiati in una qualsiasi lega metallica ma i *cerchi in lega* sono esclusivamente quelli predisposti per le ruote delle automobili.

Esistono anche espressioni in cui solo uno dei costituenti mantiene il proprio significato, mentre l'altro ne assume uno idiomatrico solo occorrendo insieme al primo: è il caso di *chiave inglese*, dove l'aggettivo che normalmente indica una nazionalità qui precisa il particolare tipo di strumento.

In generale il principio di composizionalità, di fregeiana memoria⁸, appare come requisito essenziale di molte teorie semantiche (di carattere più o meno formalizzante, prima fra tutte la Grammatica di Montague) nella versione attuale proposta nei lavori di Katz e Fodor (Katz & Fodor, 1963; Katz, 1966).

⁸Molti autori menzionano il principio di composizionalità come "principio di Frege", tuttavia è possibile attestare che non ne esistono formulazioni esplicite nei suoi scritti (cfr. K. Popper, *Unended quest. An intellectual autobiography*, Fontana, 1976, p. 198) benché lo spirito della composizionalità sia presente nei suoi ultimi lavori.

Tuttavia, per quanto a prima vista banale, la definizione di cosa sia il *non composizionale* non è ovvia e presenta non pochi problemi sia a causa della difficoltà di definizione esplicita di «cos'è una parte, cos'è un significato, quali regole di composizione o funzioni sono in gioco» (Casadei, 1996, p. 16), sia per la problematicità di attribuzione di un significato o un senso⁹ “assoluto” alle parole.

Lo stesso Frege, del resto, nell'introduzione alle sue *Grundlagen der Arithmetik*, presenta quello che in seguito prenderà il nome di *principio della contestualità*:

der Bedeutung der Wörter muss im Satzzusammenhange, nicht in ihrer Vereinzelung gefragt werden (Frege, 1884, p. xxii).

che sembra aprire precocemente il tracciato ai molti studi novecenteschi su contesto e cooccorrenza.

Firth, infatti, più di settant'anni dopo, seguirà e abbraccerà in pieno il principio della contestualità, in prospettiva distribuzionalista, con la sua famosa affermazione: «You shall know a word by the company it keeps» (Firth, 1957). Il principio che guida la lettura delle espressioni di più parole, cioè la necessità di guardare oltre il singolo costituente per interpretarne il senso, sembra innalzarsi a norma generale e la composizionalità perde ogni fondamento. In Firth (*ibid.*, p. 190) si legge:

The use of the word 'meaning' is subject to the general rule that each word when used in a new context is a new word.

Il problema della non composizionalità, quindi, travalica il confine delle espressioni qui in esame per abbracciare l'intera problematica della costruzione del senso nella lingua. Che sia a livello della parola isolata, dell'espressione idiomatica, del sintagma, della frase o del discorso, i significanti e le regole secondo cui si legano non sembrano esaurire le variabili necessarie alla comprensione del senso.

De Mauro & Voghera (1996, p. 100), in una posizione più moderata, mettono sapientemente in luce le dinamiche di interazione tra elementi linguistici e non nella creazione del senso:

Tra significato della frase e significato delle parole che la costituiscono ci pare di scorgere un rapporto di co-variabilità: ci pare indubbio il concorso dei significati delle singole parole al costituirsi del significato della frase; ma nelle parole, tra le loro molte accezioni, concorre al significato della frase quella accezione che la frase (anzitutto, ma non solo) seleziona. Similmente ci pare di cogliere un rapporto di co-variabilità tra contesto non verbale e frase nella determinazione del senso dell'enunziato: la selezione del senso avviene certamente anche in rapporto al contesto non verbale di enunciazione, ma gli elementi del contesto chiamati a fungere da selettori a loro volta sono chiamati in causa dalla forma della frase, dal suo significato e dal significato delle sue parole.

⁹Nel presente lavoro si considera la convenzione demauriana di indicare con *significato* la classe del senso al livello della *langue*, mentre con *senso* la realizzazione del significato nelle singole e particolari situazioni, quindi al livello di *parole*.

Pur essendo, quindi, la non composizionalità una caratteristica intuitivamente percepibile nell'analisi semantica delle espressioni di più parole, essa non può essere la base di una semantica formale, non solo di tali fenomeni ma della lingua stessa, in quanto presupporrebbe, in qualche modo, la visione della lingua come calcolo¹⁰ e verterebbe, inoltre, sull'autonomia linguistica del significato rispetto al contesto. Nonostante ciò, è innegabile attribuire al concetto di non composizionalità, anche se parziale, l'utilità di individuare in maniera intuitiva le espressioni il cui significato globale non ha legami stretti con le possibili accezioni¹¹ comunemente note dei suoi costituenti.

È interessante, infine, discutere di un'ultima anomalia relativa alle espressioni di più parole che De Mauro & Voghera (*ibid.*) mettono in risalto. Infatti, a differenza di ciò che avviene per le parole semplici, fortemente ambigue e profondamente soggette a selezioni di senso in base al contesto di occorrenza, le espressioni che i due autori chiamano lessemi complessi sembrano comportarsi all'opposto. Una volta che si costituisce e acquista un'identità riconosciuta nella norma d'uso, infatti, il lessema complesso «disciplina e regola un gruppo ristretto di possibili co-selezioni». Una parola come *lampo* ha, come sostantivo maschile, 7 accezioni nel GRADIT (De Mauro, 1999-2007) quali: a) fenomeno atmosferico prodotto da scariche elettriche, b) emissione luminosa intensa e di breve durata, c) flash fotografico, d) manifestazione rapida e fugace di un sentimento o di un'emozione, e) spazio di tempo, fatto o avvenimento di brevissima durata, f) persona, animale o cosa velocissima, g) rapidità fulminea, h) imbarcazione o treno veloce. Se si considera un lessema complesso in cui esso compare, come *lampo di genio*, è facile individuare l'accezione d) come l'unica che esso seleziona. In questo senso le espressioni di più parole sembrano diventare un meccanismo linguistico che aiuta, ancor prima dell'intera frase o del discorso, l'inquadramento del senso dell'enunciato.

1.3 La normalità delle espressioni di più parole

Insieme alle loro anomalie, una delle ragioni per cui le espressioni formate da più di una parola sono state a lungo relegate ai limiti di molte teorie del linguaggio è forse la percezione che esse rappresentassero delle eccezioni o comunque un fenomeno marginale rispetto alle parole semplici.

Diversi studi, tuttavia, hanno indicato che a livello quantitativo esse non costituiscono affatto dei fenomeni sparsi, isolati e marginali. Jackendoff (1997) afferma che il numero di espressioni complesse sia almeno pari a quello dei lessemi semplici memorizzati nel lessico di ciascun parlante, mentre Mel'čuk (1998) azzarda che esse superino le parole semplici in proporzione di circa 10 a 1. In ogni caso, indipenden-

¹⁰Per approfondire il concetto di lingua come calcolo si veda De Mauro, *Minisemantica dei linguaggi non verbali e delle lingue*, Laterza, 1982.

¹¹Per *accezione* si intende, qui, uno dei possibili raggruppamenti di sensi all'interno dello stesso significato.

temente da queste previsioni, è possibile ricorrere a stime più sicure che scaturiscano da dati quantitativi. Voghera (1994) illustra uno studio sui lessemi complessi compiuto sul corpus LIP (De Mauro *et al.*, 1993) in cui essi risultano 1595 su un totale di circa 16.000 lemmi. In questo caso il dato è tratto da un corpus esclusivamente di parlato, che può mostrare caratteristiche diverse rispetto a un insieme di dati linguistici che comprendano la produzione scritta; ciononostante la percentuale di lessemi complessi non è irrilevante e, anzi, costituisce un insieme tutt'altro che trascurabile.

Le espressioni di più parole, del resto, sono in un certo senso necessarie e utili alla lingua: alcuni concetti, infatti, non hanno la possibilità di essere espressi se non con un lessema complesso, come ad esempio *macchina da scrivere* o *frutto di mare*.

Molte delle congiunzioni italiane scaturiscono da processi di agglutinazione di costituenti una volta liberi, come testimoniano i seguenti esempi danteschi:

quando Virgilio incominciò: “Amore,
acceso di virtù, sempre altro accese,
pur che la fiamma sua paresse fore;
(*Purgatorio* XXII, 10-12)

saltò, forse credendo saltare uno muro, *non ostante* che 'l pastore, pian-
gendo e gridando, colle braccia e col petto dinanzi si parava.
(*Convivio*, I Trattato, XI, 10)

La cristallizzazione dell'unione dei costituenti, in questi casi, ha innescato un processo di grammaticalizzazione che ha quindi condotto alla lessicalizzazione e all'attribuzione di un nuovo significato all'espressione. In questo senso le espressioni di più di una parola grafica appaiono come «uno dei meccanismi generali e produttivi del lessico di una lingua» (Voghera, *ibid.*).

Durante la seconda metà del XX secolo, del resto, la consapevolezza che le espressioni di più parole rappresentino un fenomeno fondamentale e *normale* della lingua ha cominciato a farsi strada, in particolare in ambito anglosassone. Il principio che governa la costruzione e l'uso di tali espressioni sembra diventare, come nel caso di Firth, il principio guida della lingua. Searle ad esempio, in ambito pragmatico, arriva a formulare la nota massima «Speak idiomatically unless there is some good reason not to do so» (Searle, 1975). È solo con Sinclair (1991), tuttavia, che si arriva a una formulazione esplicita del cosiddetto *idiom principle*, che rappresenterebbe una delle due possibilità della lingua nella produzione di enunciati. Da un lato, infatti, Sinclair riconosce il principio di scelta libera (*open-choice principle*) come fondamento della grammatica, secondo cui in ogni punto della frase, l'ultima parola apre una vasta scelta di opzioni riguardo la successiva, il cui unico limite è la grammaticalità. Il discorso si struttura quindi in una successione di scelte limitate solo a livello locale dalle parole che di volta in volta occorrono. In contrapposizione a questo, il principio idiomatico guida la costruzione delle frasi in un modo differente:

The principle of idiom is that a language user has available to him or her a large number of semi-preconstructed phrases that constitute single choices, even though they might appear to be analysable into segments.

Sinclair riconosce che un tale principio può scaturire dalla ripetitività di situazioni simili, dalla tendenza al minimo sforzo o da non meglio specificate esigenze della conversazione in tempo reale. In linea con tale approccio Urzì (2009, p.1) aggiunge che «il legame privilegiato che unisce i costituenti di tali sintagmi [ciò che nel seguito sarà definito come collocazione, *nda*] finisce per stabilire tra di essi una sorta di “ponte”, creando nell’interlocutore (o nel lettore) un meccanismo di attesa che facilita la ricezione e la comprensione del messaggio linguistico». In questo senso fenomeni di idiomaticità e combinazione lessicale si configurano come attori privilegiati degli scambi comunicativi, oltre che aspetti naturali dell’interazione linguistica.

1.4 Polirematiche e collocazioni

Fino a questo punto il fenomeno delle espressioni di più parole è stato trattato come un *mare magnum* dalle molteplici e differenziate caratteristiche che mal sopporta i tentativi di categorizzazione. Tanto le anomalie, quanto le caratteristiche di normalità, sono state esposte attraverso esempi che, pur essendo molto diversi tra loro, mettevano via via in risalto gli aspetti di interesse. È tuttavia importante ricordare, come già accennato, che i legami che uniscono i costituenti delle espressioni, oltre che vari, hanno diversi gradi di intensità. Mentre *fai da te*, ad esempio, presenta un alto grado di coesione tra i costituenti, *sfatare un mito* mostra un legame più debole. Tale debolezza risiede nel fatto che esempi come quest’ultimo (o anche *occhi castani*, *pioggia torrenziale*) non costituiscono delle vere e proprie unità linguistiche, in quanto la loro occorrenza congiunta è solo preferenziale, più che necessaria. In questo insieme indistinto e variegato è utile quindi fissare, per il momento, due grandi classi di fenomeni, che sono in realtà polarizzazioni opposte di un *continuum*. In italiano si è soliti riferirsi con il nome di (*unità*) *polirematiche* alle espressioni del primo tipo, che mostrano, cioè, **necessità** di occorrenza dei costituenti al fine di veicolare uno specifico significato¹². Sono generalmente riconosciute, invece, *collocazioni* quelle

¹²Come già detto, la questione terminologica in materia non ha mai raggiunto una sintesi o prodotto una convenzione ampiamente condivisa, lasciando spazio al fiorire di innumerevoli etichette. Per quanto riguarda strettamente ciò che in questo lavoro viene definito *polirematica*, si è scelto di utilizzare la terminologia di matrice lessicografica, presente nei progetti riconducibili a De Mauro (LIP, GRADIT), Sabatini, Coletti (DISC), nonché in testi e lavori della letteratura italiana (De Mauro & Voghera, 1996; Voghera, 2004). Sempre in italiano, altri lavori utilizzano termini quali *composto sintagmatico* (Scalise, 1994), *espressione idiomatica* (Vietri, 1985; Casadei, 1996), *lessema complesso* (Voghera, 1994; De Mauro & Voghera, 1996), *espressione multiparola* (Masini, 2007), *parola sintagmatica* (Masini, 2007), pienamente assimilabili a sinonimi di *polirematica*. Altri termini come *frase fissa*, *modo di dire*, *cliché* sembrano abbracciare uno spazio più ampio che include espressioni marcate pragmaticamente, con la presenza del verbo e di dimensioni uguali o maggiori di un sintagma. In ambito anglosassone si è invece imposto l’uso contemporaneo

unioni di parole che hanno subito un “irrigidimento” nell’uso, ma non una completa cristallizzazione, e che quindi manifestano solo una forte **preferenza** di occorrenza congiunta, costituendo «un’unità fraseologica non fissa ma riconoscibile» (Tiberii, 2012, p. 3).

Mentre le polirematiche sembrano poter presentare, in generale, tutte le quattro principali anomalie elencate in precedenza, le collocazioni sembrano contraddistinte, come ipotizzato recentemente (Masini, 2009), in particolare dalla seconda, vale a dire la non sostituibilità. Il diverso comportamento è indice della diversa natura dei due fenomeni. Il nuovo significato globale che si crea con la cooccorrenza dei costituenti di una polirematica sembra imporre all’espressione una maggiore rigidità, ed essa tende a configurarsi come una vera e propria unità sul piano lessicale, semantico e sintattico. Le collocazioni, invece, manifestano un legame sintagmatico tra i propri costituenti, i quali tuttavia sono ampiamente liberi e non concorrono alla formazione di un’entità unitaria. La non sostituibilità testimonia tale legame, mostrando come per determinate parole esistano proprietà di selezione sulla base dell’uso della lingua¹³ (cfr. l’esempio citato in precedenza di *capelli castani* vs. **capelli marroni*). Si può dire che per entrambi i fenomeni sia chiara l’esistenza del collegamento tra i costituenti, ma mentre le collocazioni manifestano tale legame sulla base della frequente cooccorrenza di due o più parole, che si combinano secondo le normali regole grammaticali, le polirematiche sono il risultato di un processo di fossilizzazione lessicale e semantico (Moon, 1997). Per le polirematiche il legame, a qualsiasi livello lo si consideri, è al di sopra dei costituenti e li ingloba, mentre per le collocazioni esso rimane *tra* i componenti. È per questo che il rapporto collocazionale può essere *non reciproco*, ovvero non simmetrico per tutti i costituenti, come sottolineato da Nelson:

Blonde will only collocate with a very limited number of words - *hair* (or words that in this instance in some way relate back to *hair*, e.g. *girl*, *woman*, but *hair* will collocate with many words, e.g. *brown*, *long*, *short* and *mousy*. Thus, the strength of the bond between words is not equal (Nelson, 2000, cap. IV).

Le sfumature lungo il *continuum* che unisce polirematiche e collocazioni hanno spesso condotto gli studiosi a raggruppare insieme i fenomeni. Le polirematiche o le frasi completamente cristallizzate, a causa della loro rigidità semantica e sintattica sono spesso apparse come unità ideali e di conseguenza le collocazioni sono diventate un loro sottoinsieme. Secondo Mel’čuk, infatti «collocations - no matter how one understands them - are a subclass of what are known as set phrases» (Mel’čuk, 1998,

predominante e generalizzato di *multiword expression*, benché in letteratura siano presenti anche termini quali *idiom* (Malkiel, 1959) o *multi-word item* (Nelson, 2000).

¹³Cfr. la definizione di van Roey: «[collocation is] that linguistic phenomenon whereby a given vocabulary item prefers the company of another item rather than its ‘synonyms’ because of constraints which are not on the level of syntax or conceptual meaning but on that of usage» (van Roey, 1990, p. 46).

p. 23). Van der Wouden (1997), al contrario, concede alle collocazioni uno status preminente, suggerendo la nozione di *collocabilità* piuttosto che di *idiomaticità*. Per van der Wouden con il termine collocazione ci si riferisce a tutti i tipi di combinazioni lessicali fisse ed in particolare anche alle polirematiche, che diventerebbero «those collocations with a non-compositional or opaque semantics»¹⁴ (Van der Wouden, 1997, p. 9).

Più in generale è possibile identificare polirematiche e collocazioni come polarizzazioni opposte di un segmento continuo. Howarth (1996, p.32-33), nel suo modello di *continuum* lessicale segue questa strada, proponendo quattro gruppi categoriali che rappresentano i gradi di separazione tra combinazioni relativamente libere e fisse:

Free collocation	<i>blow a trumpet</i>	‘suonare la tromba’
Restricted collocation	<i>blow a fuse</i>	‘far saltare un fusibile’, o (idiom.) ‘arrabbiarsi’
Figurative idiom	<i>blow your own trumpet</i>	‘vantarsi in maniera eccessiva’
Pure idiom	<i>blow the gaff</i>	‘rivelare una verità nascosta’

I criteri seguiti nella costruzione di una tale suddivisione risultano sia lessicali (prime due categorie) che semantici (ultime due). I gruppi, inoltre, presentano anch’essi confini sfumati, data la problematicità di identificazione di *syntax restrictedness*, *figurativeness* o dell’opacità semantica invocata da Howarth per l’ultima categoria.

1.5 Approcci classici ai fenomeni idiomatico-collocativi

Nella storia degli studi linguistici il fenomeno delle “espressioni di più parole” è comparso più volte in varie teorie sulla lingua, subendo diverse interpretazioni e analisi, ognuna delle quali ha aggiunto o modificato definizioni, ambito e strumenti di analisi di tali manifestazioni linguistiche.

Già nel XIX secolo, è possibile rintracciare i primi sentori dell’importanza di entità più grandi delle singole parole, ma non abbastanza da abbracciare lo status di frase. Prendergast (1864), in un suo lavoro in ambito pedagogico, nota come nell’apprendimento della lingua i bambini imparino «not just words, but ‘chunks’ of language» (Nelson, 2000), utilizzandoli fluentemente nei loro discorsi.

Appena più tardi Paul (1880) pone l’attenzione su alcune combinazioni fisse del tedesco che esibiscono un comportamento peculiare, non ammettendo sostituzioni di costituenti, modificazioni e il cui significato non è ricostruibile dai componenti. Per tale ragione queste espressioni sembrano equivalere a concetti unitari.

¹⁴Si ricorda, tuttavia, che la non composizionalità non è requisito essenziale di un’espressione complessa quale una polirematica. Si pensi ad es. a *presidente del consiglio* o *scuola elementare*, entrambe espressioni iponime della testa, e quindi composizionali, ma che identificano unità semantiche ben precise e per questo riconoscibili come polirematiche.

Sweet (1891) attesta l'esistenza, in inglese, di *special sentences, or idioms* che, sebbene costruiti regolarmente nella forma, al pari delle frasi standard, sono irregolari nel significato, poiché «the meaning of the whole cannot be inferred from the meaning of its elements» (*ibid.*, p. 156). La non composizionalità, quindi, è ancora una volta la più evidente peculiarità indicativa delle “espressioni irregolari”.

Poco più tardi, Bréal (1904) pone l'attenzione sui *groupes articulés* e cioè locuzioni fisse ricorrenti divenute usuali di cui «nous ne percevons plus que la formule». La perdita percettiva dei singoli elementi è favorita, secondo Bréal, dal ricorso a strutture sintattiche disusate o anche la conservazione di significati lessicali obsoleti¹⁵ (Casadei, 1996, p. 33).

Nel XX secolo l'attenzione su lessemi complessi, espressioni idiomatiche, legami lessicali si intensifica, seppure non si riesca a sviluppare un trattamento del fenomeno che risulti esaustivo e strutturato. Le unità polirematiche rimangono ai margini delle molte teorie linguistiche sviluppate dalle principali scuole di pensiero poiché viste come anomalie o eccezioni che mal si adattano a trattazioni sistematiche. Il concetto di collocazione, invece, si svilupperà solo nella seconda metà del secolo, contribuendo a stimolare la riflessione sui legami tra le parole, specialmente sul piano lessicale. Nel seguito è esposta una breve trattazione storica dei principali approcci al fenomeno sviluppati nel secolo scorso¹⁶.

1.5.1 La visione strutturalista

In ambito strutturalista l'attenzione alle espressioni di più di una parola converge col tema centrale dell'identificazione delle parti, delle unità linguistiche e dei rapporti che esse intrattengono con gli altri componenti. Masini (2007, p. 8) osserva come «la strutturazione in livelli di analisi e l'identificazione delle unità linguistiche siano uno dei lasciti più importanti dello strutturalismo» ed è proprio per questo che i fenomeni di confine tra domini, espressioni di più parole *in primis*, risultano profondamente stimolanti per la riflessione linguistica. In ogni caso, lo studio del fenomeno idiomatico da parte degli strutturalisti si pone storicamente come la base per ogni altro approccio successivo.

Riferimenti alle espressioni idiomatiche si possono rintracciare in Saussure, che nel *Cours* accenna a combinazioni di parole appartenenti all'ambito della *langue*:

On rencontre d'abord un grand nombre d'expressions qui appartiennent à la langue; ce sont les locutions toutes faites, auxquelles l'usage interdit de rien changer, même si l'on peut y distinguer, à la reflexion, des parties significatives [...]. Ces tours ne peuvent pas être improvisés, ils sont fournis par la tradition (Saussure, 1922, p. 172).

¹⁵Si pensi all'espressione italiana *chiedere venia* dove *venia* è l'equivalente obsoleto di *perdono*.

¹⁶La trattazione va intesa come una carrellata sulle principali idee e i principali autori che hanno studiato i fenomeni delle espressioni di più parole. Per una visione dettagliata, comprensiva di valutazioni critiche, si rimanda ai lavori di Casadei (1996, pp. 27-80) e Masini (2007, pp. 7-40).

Una tale affermazione sembra contenere *in nuce* l'*idiom principle* di Sinclair, in quanto il parlante, che gestisce di norma la sua particolare interazione linguistica con frasi libere sul piano della *parole*, può fare ricorso, quando serve, a locuzioni precostruite e messe a sua disposizione dalla tradizione. Appena più avanti si precisa, infatti, la difficoltà nello stabilire quanto un certo sintagma¹⁷ sia considerato un fatto di *langue* o di *parole* poiché sia l'uno che l'altro fattore possono contribuire alla sua formazione in proporzioni difficilmente determinabili, allo stesso modo di come l'*idiom* e l'*open-choice principles* regolano la produzione linguistica.

Un altro spunto alla riflessione sulle espressioni complesse compare anche nella trattazione del fenomeno dell'agglutinazione, dove Saussure esplicita ciò a cui si è già accennato riguardo a delessicalizzazione e accorpamento di sequenze ripetute generalmente in blocco, che tendono a diventare parole semplici:

[...] quand un concept composé est exprimé par une suite d'unités significatives très usuelle, l'ésprit, prenant pour ainsi dire le chemin de traverse, renonce à l'analyse et applique le concept en bloc sur le groupe de signes qui devient alors une unité simple (Saussure, *ibid.*, p. 243).

Il piano lessicale e quello grammaticale sono quindi in contatto e lessico e grammatica figurano come due poli del *continuum* linguistico, due "correnti" che spingono e condizionano la produzione. Il lessico è la polarità dell'"immotivato", ovvero dell'arbitrario e inanalizzabile; la grammatica è il fulcro della motivazione e delle regole di costruzione¹⁸ (Saussure, *ibid.*, p. 183).

Inoltre la frequente sostituibilità (sia a livello intra- che interlinguistico) tra forme semplici e costruzioni perifrastiche mette in luce quanto le unità lessicali abbiano una similarità funzionale con le unità grammaticali. Infatti è centrale, nella visione strutturalista, l'idea dell'espressione idiomatizzata come unità commutabile con parole semplici in quanto espressione di concetti unitari.

Sechehaye (1921) sottolinea che le locuzioni rappresentano un elemento problematico nella grammatica della lingua proprio perché alienate dal sistema morfologico e semantico. L'unitarietà della struttura (e il carattere funzionale di unità) annulla l'identità dei singoli costituenti e quindi riporta l'espressione allo stesso livello della parola semplice, in quanto è solo quest'ultima a svolgere primariamente il ruolo di veicolo di singoli significati. Le espressioni idiomatiche, infatti, intrattengono rapporti associativi con le parole semplici, come sottolineato da Greimas (1960): *pomme de terre* si oppone paradigmaticamente ad esempio a *prune* o *betterave* e

¹⁷Per Saussure il sintagma «se compose [...] toujours de deux ou plusieurs unités consécutives» e inoltre «s'applique non seulement aux mots, mais aux groupes de mots, aux unités complexes de toute dimension et de toute espèce (mots composés, dérivés, membres de phrase, phrase entières» (Saussure, 1922, pp. 170-172).

¹⁸«Non que "lexique" et "arbitraire" d'une part, "grammaire" et "motivation relative" de l'autre, soient toujours synonymes; mais il y a quelque chose de commun dans le principe» (Saussure, 1922, p. 183).

per questo essa deve essere considerata *parola*, unità lessicale. Gremais, inoltre, concede alla semantica di operare una distizione tra le espressioni che conservano ancora un qualche legame con i costituenti (come *ferro da stiro*) da quelle completamente indipendenti (come nel caso di *luna di miele*).

I diversi studi cominciano quindi ad accrescere la terminologia. Benveniste (1966) conia il termine *sinapsi* per indicare unioni di due lessemi mediante una preposizione (diversamente da quanto avviene per i composti), come *chiaro di luna*: esse si distinguono dai sintagmi liberi in base al referente (entità unitaria per le prime, insieme di concetti distinti per i secondi). Appena più tardi Martinet suggerisce l'uso di *sintema* per indicare «les unités linguistiques dont le comportement syntaxique est strictement identique à celui des monèmes avec lesquels ils commutent, mais qui peuvent être conçue comme sémantiquement analysable» (Martinet, 1967, p. 6). Egli crea, quindi, una categoria *ad hoc* per inglobare tutte le unità che non possono rientrare in morfemi o monemi¹⁹, in quanto aventi una struttura e non essendo elementi che Saussure chiamerebbe “immotivati”.

Il fatto che le espressioni più o meno cristallizzate si oppongano paradigmaticamente e sintagmaticamente ad altre entità del discorso, induce Coseriu (1966), in un'ottica totalmente diversa da quella di Martinet, a negare loro una propria identità, riconducendole allo status di frase, sintagma o parola a seconda della possibile commutabilità che esse manifestano con questi ultimi.

Esiste tuttavia una gradazione di coesione interna dei costituenti che pone un limite all'assoluta interscambiabilità delle espressioni idiomatiche con le parole. Bally (1951) è forse il primo a riconoscere che oltre alle espressioni fisse, da lui definite *unités phraséologiques*, esistono *séries phraséologiques* assimilabili in parte all'odierno concetto di collocazione. Mentre per i primi, infatti, il significato dell'espressione è globale e indipendente dai costituenti, i secondi rappresentano modi di dire abituali o soltanto ricorrenti. Il lavoro di Bally, straordinariamente esteso e ragionato, risulta il primo nell'ambito di un nuovo filone di studi che darà origine alla fraseologia.

La particolare indipendenza del significato rispetto ai costituenti di molte espressioni induce a riflettere meglio su ciò che si intende per idiomatico all'interno di una lingua²⁰. A metà degli anni cinquanta, infatti, Hockett (1956) formula un nuovo concetto di espressione idiomatica, definendola come una forma grammaticale il cui significato non è deducibile dalla propria struttura. Quest'ottica, tuttavia, ha una conseguenza importante: anche l'intera categoria dei morfemi viene inglobata nelle espressioni idiomatiche, che diventano quindi un elemento fondamentale della struttura di una lingua. Ricollegandosi all'intuizione di Bréal, secondo cui l'obsolescenza

¹⁹Rispettivamente unità minime sul piano della forma e del significato.

²⁰Nella tradizione linguistica sono rintracciabili due accezioni di *idiomatico*: la prima è l'idiomaticità *interlinguistica*, in cui un'espressione risulta idiomatica perché confrontata con analoghe in altre lingue; la seconda è l'idiomaticità *intra-linguistica* per cui un'espressione è definita idiomatica mediante il confronto con forme o strutture di quella stessa lingua. Tuttavia solo la seconda sopravvive come principale modo d'intendere l'idiomatico. Per approfondire, cfr. Casadei (1996, pp. 27-28).

del lessico è una caratteristica peculiare in molte espressioni idiomatiche, Hockett ipotizza che le espressioni abbiano maggiore possibilità di divenire idiomatiche se create a partire da *pattern* in regresso. Inoltre, il campo dell'idiom invade anche tutte le forme semplici poiché un altro modo di identificazione delle espressioni idiomatiche è considerare se, a seconda del contesto d'uso, esse hanno un diverso referente. In questo modo termini anaforici, nomi propri, pronomi, diventano unità idiomatiche, al pari delle espressioni di più parole. Questa visione tuttavia, pur rendendo «la nozione di idiomtico pervasiva e rilevante, oltre che del tutto non anomala», fa sì che essa diventi «troppo ampia e di fatto inutilizzabile» (Casadei, 1996, p. 39).

L'intuizione di Hockett di considerare l'idiomatico e le espressioni idiomatiche non più come un'anomalia viene ripresa nei successivi filoni strutturalisti americani nei modelli tagmemici e stratificazionali. L'idea di base che permette una tale visione è l'indipendenza dei livelli linguistici: se infatti ogni strato ha completa autonomia rispetto agli altri, il livello semantico può costruire le proprie unità indipendentemente dal livello morfologico o sintattico, non avendo quindi, come limite, i confini delle parole semplici. Per Pike (1967) le espressioni idiomatiche sono ipermorfemi²¹ perfettamente regolari a livello sintattico ma con un significato non compositivo. Grazie all'autonomia dei livelli è possibile che nelle espressioni idiomatiche siano conservati i morfemi (e quindi la possibilità di flessione), ma non i sememi, ovvero i tratti semantici dei costituenti (Makkai, 1972).

A Makkai si deve infine anche la nota suddivisione delle espressioni idiomatiche in due principali categorie, ovvero gli *idioms of encoding* e gli *idioms of decoding*. I primi abbracciano tutte le espressioni che nella loro struttura utilizzano costruzioni particolari, come l'uso peculiare di una preposizione. In italiano è questo il caso ad esempio di *da te* nella frase *vengo da te* dove la preposizione *da* indica moto a luogo. L'encoding quindi rappresenta una scelta particolare e irrinunciabile di costruzione che però si esplica solo in una o poche espressioni mediante l'uso di lessico o *pattern* in genere adoperati in modo diverso nella sintassi standard. Gli *idioms of decoding*, invece, raggruppano le espressioni che pongono un problema di decodifica semantica in quanto interpretabili sia idiomaticamente che secondo il loro senso letterale. Il problema di decodifica, quindi, fa riemergere un certo grado di anomalia, anche se, per la prima volta, puramente nell'ottica della comprensione.

Il quadro eterogeneo che scaturisce dalla varietà di visioni degli studiosi strutturalisti testimonia quanto il fenomeno dell'idiomaticità sia presente nella riflessione linguistica moderna fin dal principio. Ben lungi dal fornire risposte definitive ed esaustive, la riflessione strutturalista ha il merito di aver attestato l'esistenza di unità intermedie tra parola e frase e di aver sottolineato la problematicità del trattamento delle loro caratteristiche sintattiche e semantiche a cavallo dei diversi livelli

²¹In tagmemica, ogni tagmema indica l'insieme degli elementi che possono occupare una certa posizione sintagmatica. Questi possono essere unità semplici o complesse, dette in quest'ultimo caso *ipermorfemi*.

di analisi considerati.

1.5.2 La visione generativista

In generale in ambito generativista l'attenzione è posta principalmente su espressioni idiomatiche verbali, che vengono considerate strutture ben formate grammaticalmente ma con anomalie semantiche. Tali anomalie, insieme alle restrizioni sulle trasformazioni ad esse applicabili, sono il principale nodo di riflessione poiché si configurano come i principali ostacoli all'inclusione delle espressioni nella teoria standard chomskyana (Vietri, 1985), e si tenta spesso, perciò, «di risolvere nel modo più economico il problema idiomatico piuttosto che comprenderlo» (Casadei, 1996, p.43), per dedicarsi piuttosto ad approfondire lo studio della *core grammar* (Masini, 2007). Se il lessico è solo «an appendix of the grammar» (Bloomfield, 1967), ad esso è relegato l'unico scopo di fornire passivamente una lista di parole o di espressioni a cui attingere per riempire gli *slot* della struttura sintagmatica, ovvero i nodi terminali (X^0). Masini (2007) sottolinea che l'anomalia idiomatica prorompe proprio a questo stadio del processo: se infatti i nodi X^0 , in quanto terminali, presuppongono *item* semplici, le espressioni idiomatiche che vanno ad occuparli sono obbligate a privarsi della propria struttura interna, che invece sembra sopravvivere nei casi in cui varie trasformazioni sono possibili.

Nei primi lavori sull'argomento (Katz & Postal, 1963; Weinrich, 1966, 1969) si ipotizza l'inclusione delle espressioni idiomatiche nel dizionario lessicale. Mentre in Katz & Postal (1963) ogni espressione, da inserire nella struttura come un tutt'uno (inserzione globale unitaria), va ricondotta a una sola categoria (verbo, nome, ecc.) e sono escluse le espressioni malformate, in Weinrich (1969) si suggerisce di includere queste ultime nel dizionario di lemmi semplici e di segnalare i tratti di anomalia trasformativa accanto all'espressione di riferimento. Weinrich è inoltre promotore della cosiddetta *tesi dell'azzeramento semantico* secondo cui non c'è rapporto tra il significato idiomatico e quello letterale di un'espressione, per cui l'idiom diventa semplicemente un omonimo dell'espressione a lettura compositiva. Una apposita *idiom comparison rule* avrà il compito di attivare la lettura idiomatica rispetto a quella compositiva standard.

Sul tema della compositività Weinrich sottolinea come non esista nemmeno isomorfismo tra i costituenti dell'espressione idiomatica e le parole di una possibile parafrasi: per i morfemi di *tirare le cuoia* non possiamo identificare una corrispondenza biunivoca con quelli (meno numerosi) di *morire*. Anche per Chafe (1968) i costituenti dell'espressione sono assenti a livello semantico e per questo non possono essere modificati. Fraser (1970) suggerisce che i singoli costituenti abbiano un'informazione sintattica che permetta loro di possedere una struttura, tuttavia l'informazione semantica sarebbe totalmente assente. Se, quindi, i componenti non partecipano singolarmente alla creazione del senso dell'espressione, per Fraser è inevitabile che vengano bloccate anche le trasformazioni sintattiche operabili di nor-

ma sulle parole autonome: è la semantica, quindi, a governare alcuni aspetti della sintassi.

Sebbene «anacronistico da un punto di vista teorico» (Masini, 2007, p. 31), uno studio di Schenk (1995) riprenderà le idee di Fraser, ipotizzando che esistano trasformazioni sintattiche strettamente legate alla semantica dei costituenti, come topicalizzazione, modificazione, pronominalizzazione, che perciò non possono operare su espressioni idiomatiche non composizionali. Trasformazioni come *raising* o *verb second* (V2), invece, agendo ad un livello puramente grammaticale, possono interessare qualsiasi frase, indipendentemente dal contenuto semantico dei componenti.

Il rapporto tra semantica e sintassi e l'influenza della prima sulla seconda rimangono punti centrali nell'evoluzione delle riflessioni sugli *idioms* e sulle trasformazioni bloccate. Poco dopo Fraser, Newmeyer (1972, 1974) suggerisce che la soluzione alle restrizioni di trasformazione stia nel considerare il rapporto tra i significati composizionale e idiomatico di ogni espressione. Infatti, le regole cicliche che governano le trasformazioni sembrano possibili solo quando applicabili sia al significato letterale che alla parafrasi di quello idiomatico: se *tirare le cuoia* non è passivizzabile, dipende dal fatto che *morire* non lo è. Al contrario *rompere il ghiaccio* ammette il passivo perché lo fa anche la sua interpretazione *dissipare l'imbarazzo*²². In seguito Ruwet (1983) riprenderà il rapporto tra trasformazioni e composizionalità del significato, affermando che la libertà di un'espressione idiomatica dipende dal suo grado di decomponibilità. È possibile infatti, invece di considerare l'intera parafrasi, associare ad ognuno dei costituenti idiomatici un componente di significato: in *rompere il ghiaccio*, *rompere* significa *dissipare* e il *ghiaccio* è *l'imbarazzo*. Tutte le espressioni decomponibili in questo modo possono essere modificate.

Chomsky stesso si sofferma sul tema delle espressioni idiomatiche, citandole come «prova schiacciante» della sua concezione di *struttura profonda* e *struttura superficiale* (Chomsky, 1980). Le espressioni idiomatiche, infatti, testimoniano l'asimmetria dei due piani poiché sembrano essere presenti o a entrambi i livelli, o nella struttura profonda, ma mai solo al livello di struttura superficiale (non consentendo, ad esempio, la passivizzazione). Poiché, inoltre, le espressioni sono ben formate e spesso hanno una lettura composizionale, la sintassi le genera allo stesso modo delle frasi libere, ma solo in seguito una *idiom rule* assegna loro il valore idiomatico.

Dalla trattazione chomskyana restano escluse le molte eccezioni considerate irrilevanti, ma in realtà cospicue, come le espressioni malformate o quelle idiomatiche solo al passivo (Wasow *et al.*, 1983; Brame, 1984a,b). In particolare sia Brame (1984a) che Ruwet (1983) forniscono controesempi di espressioni unicamente al passivo, come *il dado è tratto*, falsificando la prova schiacciante di Chomsky.

Sul versante sintattico sono quindi ampie le critiche alle idee nate in ambito generativo. L'ipotesi di un tratto *+idiom* da assegnare ad alcuni sintagmi non sembra supportata da meccanismi che possano essere inclusi nella teoria standard e rimane

²²Rimane problematica, in questa visione, la questione della parafrasi. Essa infatti non è sempre univoca o non sempre possibile.

dubbia la ragione per cui i parlanti abbiano piena competenza nello scegliere quali trasformazioni sono permesse e quali vietate, data un'espressione. Il punto focale delle critiche è inoltre la concezione dell'idiomatico come un fenomeno periferico e minoritario della produzione linguistica, tipica della visione generativa, che condurrà alla rivalutazione del lessico e dell'uso linguistico e all'importanza dei riferimenti al contesto di molte espressioni.

1.5.3 Critiche al modello generativo e interpretazioni pragmatiche

A partire dagli anni settanta alcuni studi cominciano a porsi in maniera critica nei riguardi del paradigma generativo proprio grazie ad analisi sull'idiomatico. Gross (1981, 1984), pur sviluppando i suoi lavori all'interno della teoria formale del Lessico-Grammatica, attesta che il comportamento delle espressioni idiomatiche verbali non differisce molto da quello dei sintagmi verbali "standard". Non è possibile creare un apparato sintattico generale in cui le entità rispondano a regole ben precise poiché ogni verbo, infatti, sarebbe peculiare, avendo un comportamento argomentale diverso da tutti gli altri e quindi le frasi "standard" non risulterebbero così dissimili dalle espressioni idiomatiche. In quest'ottica frasi libere e frasi fisse diventerebbero poli di un *continuum* (Masini, 2007).

Inoltre, aggiungendosi al livello semantico e sintattico, il livello pragmatico viene a configurarsi come una nuova spinta all'interpretazione delle espressioni idiomatiche. Inizia a nascere, infatti, l'ipotesi che gli idioms possano avere legami specifici con il loro contesto d'uso e che non tutto il significato o le loro anomalie dipendano intrinsecamente dalla loro natura. Nonostante siano ancora collocabili in ambito generativo, Fillmore *et al.* (1988) sono i primi a riconoscere che esistano *idioms with a pragmatic point* come nel caso di *How do you do?* Il carattere dello studio di Fillmore, Kay e Connor è più che altro tipologico, fissando quattro assi di opposizione che vedono il contrapporsi, oltre ai già citati idioms con valore pragmatico e non, di *encoding* vs *decoding*, espressioni sintatticamente corrette vs. malformate e infine *substantive* vs. *formal idioms*. Questi ultimi rappresentano pattern abituali in grado di generare un numero infinito di frasi, il cui significato non è deducibile dalla struttura o che hanno una particolare valenza pragmatica. Le idee di Fillmore *et al.* gettano le basi del filone della *Construction Grammar*, che rivaluta il concetto di "costruzione" rispetto al sistema formale chomskyano, in cui la struttura profonda risponde a un nucleo minimo di regole e la fenomenologia di variazioni delle frasi è ricondotta meramente a trasformazioni. L'idea che esistano costruzioni preimposte (come l'inglese *the Xer, the Yer*) sottolinea l'impossibilità di trascurare unità più grandi dei costituenti minimi e privi di struttura e la necessità di costruire un modello in cui sintassi, semantica, lessico e pragmatica siano altamente correlate.

Nell'ambito della semantica compositiva, Nunberg (1978), il cui lavoro era stato la base degli studi compositivi di Ruwet, suggerisce l'importanza delle co-

noscenze condivise (*normal beliefs*) nell'interpretazione delle espressioni idiomatiche (e conseguentemente nel loro largo uso). La scelta dei costituenti non è casuale, ma dipende da sfumature semantiche che essi condividono col senso globale che l'espressione può assumere. *Fare i salti mortali* implica infatti movimento, tratto condiviso dall'interpretazione dell'espressione "aver fatto di tutto per raggiungere un determinato scopo". I costituenti, quindi, ad un qualche livello, contribuiscono al senso generale. La fissità dell'espressione, inoltre, dipende dalla sua non decomponibilità. In *The Pragmatics of Reference* sono ipotizzate tre classi di espressioni: *normalmente decomponibili*, *anormalmente decomponibili* e *non decomponibili*. Mentre le prime due comprendono espressioni i cui singoli componenti possono essere associati a singoli concetti della parafrasi interpretativa (in maniera diretta per la prima, con metafore condivise per la seconda), la terza raggruppa le espressioni per cui ciò non è possibile e che, come ribadisce Ruwet, risultano fisse.

Uno studio sull'idiomatico che si colloca pienamente in ambito pragmatico è quello di Searle (1975), secondo il quale le espressioni che hanno un significato idiomatico accostano al *sentence meaning* (il significato letterale), un *utterance meaning*. Quest'ultimo annulla il primo e, in base al contesto d'uso, vi si sostituisce in un particolare senso, a differenza degli atti indiretti in cui l'*utterance meaning* viene semplicemente ad aggiungersi al significato letterale²³. Searle, tuttavia, espande la nozione di idiomatico a tutto ciò che è convenzionale nella lingua, e la sua già citata affermazione che suggerisce di parlare idiomaticamente in tutti i contesti non è in realtà riferita unicamente agli idioms. *Perché non chiudi la finestra?* è una frase idiomatica poiché è un modo convenzionale di chiedere di compiere l'azione, nonostante il *sentence meaning* esprima semplicemente l'interrogativo sul perché l'ascoltatore non lo stia facendo. In questa visione, se il parlante non usufruisce delle formule o dei pattern convenzionalmente utilizzati, egli lo fa per un preciso motivo, marcando in modo specifico la frase.

Mel'čuk (1998) ritorna sull'argomento, riportando l'attenzione nello specifico alle espressioni fisse, con la nozione di *pragmatemi*. Essi si formano quando sia il significato che il significante di un'espressione non sono liberamente e regolarmente costruiti²⁴ o quando solo il significato non lo è. La forza pragmatica del contesto d'uso gioca quindi un ruolo fondamentale:

For instance, one sees on a restaurant sign *Caesar Salad: All you can eat*; its counterpart in French sounds *Salade César à volonté* [...]. It would be semantically and syntactically correct to say in French *Salad César*:

²³Greciano (1983) contesta il processo di falsificazione del senso letterale al fine di sostituirvi quello idiomatico, poiché il primo soggiace all'espressione anche nella lettura non compositiva. Tra i due significati, quindi, esisterebbe un rapporto di omonimia sullo stesso significante.

²⁴C'è una *unrestricted construction* se i componenti vengono scelti in base a regole di selezione scelte arbitrariamente. Ciò implica che non si sia dipendenza dal contesto o dalla situazione nell'espressione del significato e che si possa liberamente scegliere la combinazione di lessemi sul piano del significante. La *regularity*, invece, implica composizionalità del significato e sintassi standard sul piano del significante (Mel'čuk, 1998, p. 3).

Tout ce que vous pouvez manger; however, this expression smacks of a calque: this is not the way the Frenchmen say it (Mel'čuk, 1998, p. 5).

Tutte le espressioni abituali e convenzionali come i saluti, le frasi da lettera, proverbi e simili sono pragmatemi, infatti «even if semantically and syntactically they are 100 percent compositional, [...] they are non-compositional pragmatically» (*ibid.*, p. 6).

1.5.4 Livello del lessico e collocazioni

Alcuni studi proposti nell'ambito dello strutturalismo europeo hanno indagato più a fondo i legami atti a solidificare l'unione tra parole, anche non mirando allo studio di unità di significato chiaramente identificabili. Si nota, infatti, che la semantica è in grado di poter generare restrizioni sul lessico (e quindi sulla scelta di accostamento delle parole), oltre che sulla sintassi.

Già Porzig (1934) formalizza il concetto di *Wesenhafte Bedeutungsbeziehungen*, ovvero *rapporti semantici essenziali* per spiegare il regolare accoppiamento di parole del tipo *gatto-miagolare*, in cui la scelta del lessico è fortemente ristretta collocazionalmente. Coseriu (1967) ritorna sull'argomento, analizzando ciò che lui definisce *Lexikalische Solidaritäten* in modo più rigoroso, focalizzandosi fortemente sulle implicazioni semantiche. Esse comprendono tutte le coppie di lessemi in cui un tratto semantico distintivo di uno dei due è compreso nell'altro e non è quindi possibile esprimere l'uno senza che sia presente l'altro nel suo significato²⁵: *calciare-piede*, *abbaiare-cane* sono esempi di questo tipo. Esiste tuttavia una differenza: in *abbaiare-cane* entrambi i lessemi intrattengono associazioni paradigmatiche con altri lessemi per cui, al variare di uno (ad esempio con la sostituzione di *cane* con *gatto*) è possibile sceglierne un altro omologo nel paradigma a disposizione (*miagolare*): tali solidarietà lessicali vengono definite da Coseriu *multilaterali*. Nel caso, invece di *calciare-piede* non esiste un paradigma analogo in quanto *calciare* non ha corrispondenti non collegati a *piede* e si è, perciò, in presenza di una *solidarietà unilaterale*.

Dieci anni prima di Coseriu, Firth (1957) aveva già pubblicato un volume che raccoglieva diversi saggi di linguistica scritti tra il 1934 e il 1951. Nel più famoso della raccolta, *Modes of Meaning*, viene per la prima volta espresso il concetto di collocazione come uno dei modi di significazione, in un'ottica opposta a quella di Porzig e Coseriu, proponendo stavolta un lessico in grado di condizionare la semantica. Secondo Firth, il significato di una parola è costruito su più livelli, e un contributo fondamentale al suo inquadramento è dato dal contesto sintagmati-

²⁵Masini (2007, p. 18), nel caso della coppia *mordere-denti* ricorda come anche Lyons (1977, p.262) accenni a *denti* come “encapsulated” in *mordere*, o Pustejovsky (1995, p. 63) che parla di *argomento ombra*.

co usuale di occorrenza²⁶. La ricorrente vicinanza di parole sull'asse sintagmatico conduce, infatti, a una specializzazione del senso dei componenti di un'espressione. In *caffè forte* l'aggettivo assume un'accezione specifica dovuta unicamente alla presenza di *caffè*. Il legame lessicale che viene a instaurarsi, quindi, fa sì che uno dei modi di definizione del significato di una parola sia l'insieme delle parole con cui essa solitamente compare, in un'ottica tipicamente distribuzionalista, come fa notare Masini (2007, p. 16). Poco più tardi Halliday (1966) riprende il lavoro di Firth, ponendo l'accento sulla necessità dell'entrata del lessico tra i livelli dell'analisi linguistica. Semantica e sintassi, infatti, non sembrano riuscire a gestire e governare *in toto* i contesti collocazionali in termini di linguistica formale poiché gli elementi che dovrebbero formare classi sul piano lessicale trascendono le relazioni cui si è soliti fare riferimento in termini di grammatica. L'idea è che esista una *lexicalness* che non sia opposta, bensì complementare alla *grammaticalness*²⁷. Halliday cita il noto esempio di *strong* e *powerful*, i quali agiscono da sinonimi o comunque intrattengono un rapporto paradigmatico in molte situazioni, come ad esempio in *a strong* o *a powerful argument*. Tuttavia non sempre la loro sostituzione è valida: è possibile *a strong argument* o *a strong tea* ma non **a strong car*; è corretto *a powerful car* ma non *a powerful tea*. Questo mostra che «the paradigmatic relation of *strong* to *powerful* is not a constant but depends on the syntagmatic relation into which each enters» (Halliday, 1966, p. 150). In termini grammaticali, la possibile spiegazione delle restrizioni collocazionali in base a tratti distintivi dei sostantivi determinati non sembra reggere. *Strong* può essere associato tanto ad astratti (*argument*) quanto a concreti (*table*) e anche raffinando i criteri di selezione, *tea* e *whisky* possono figurare entrambi come bevande o liquidi, tuttavia insieme a *strong tea*, in inglese si ha *a powerful whisky*. Halliday nota inoltre che *strong*, *strongly*, *strength* and *strengthened* hanno lo stesso comportamento collocazionale, al pari di *power*, *powerful*, *powerfully* e possono quindi essere considerati come un unico *item* a livello lessicale. La struttura di analisi di tale livello risulta molto semplice: essa si esaurisce in un grande numero di insiemi (disgiunti, o intersecantisi) che definiscono contesti collocazionali (l'insieme dei contesti possibili data una parola) e *item* lessicali e che, quindi, «are mutually defining».

Halliday è anche il primo a identificare chiaramente la necessità di una prospettiva statistica nello studio delle collocazioni, in termini di probabilità di cooccorrenza sulla base della frequenza degli *item* in un testo:

Collocation is the syntagmatic association of lexical items, *quantifiable*, *textually*, as the probability that there will occur at *n* removes (a distance of *n* lexical items) from an item *x*, the items *a*, *b*, *c*... Any given item thus

²⁶Firth tiene a precisare che il contesto a cui si riferisce è quello puramente linguistico, dato che «meaning by collocation is not at all the same thing as contextual meaning, which is the functional relation of the sentence to the processes of a context of situation in the context of culture» (Firth, 1957, p. 195).

²⁷«If therefore one speaks of a lexical level, there is no question of asserting the 'independence' of such a level, whatever this might mean» (Halliday, 1966, p. 152).

enters into a range of collocation, the items with which it is collocated being ranged *from more to less probable* (Halliday, 1961, p. 276, c.v.o mio).

Di questo avviso è anche Sinclair, che formalizza le nozioni di *node* e *span*, rispettivamente *item* di cui si vogliono analizzare le collocazioni e il numero di unità lessicali antecedenti e successive all'*item* in esame (Sinclair, 1966, p.415). Tali concetti risultano basilari per analisi collocazionali su grandi quantità di testi trattabili computazionalmente, le quali apriranno la strada all'utilizzo dei *corpora* in studi di questo tipo. In questa prospettiva, che Gledhill (2000) definisce 'statistico-testuale', le collocazioni vengono riconosciute in base a criteri squisitamente statistici, che coinvolgono la frequenza con cui i collocati si presentano nello *span* del *node* considerato, rispetto a quella con cui essi occorrono in collocazioni con altri *nodes*. Il calcolo di tali frequenze porterà Sinclair a definire le importanti nozioni di *upward* e *downward collocations* (Sinclair, 1991). Esistono infatti parole che formano generalmente collocazioni con altre che risultano molto più frequenti: è il caso, in inglese, di *back*, che di solito ricorre con *down*, *at*, *from*, *on*, preposizioni tutte più frequenti dell'avverbio; in riferimento a *back* queste collocazioni vengono definite *upward*. Allo stesso tempo, *back* compare in collocazioni con verbi come *bring* o *arrive*, che hanno una frequenza minore: si parla in questo caso di collocazioni *downward*. Sinclair si spinge a considerare la generale differenza nella natura dei due tipi di collocazioni. Mentre le *upward* sono in genere pattern più deboli in termini statistici e sono rappresentate, in larga parte, da associazioni sintattiche (preposizioni, congiunzioni, avverbi), le *downward collocations* aggiungono un tratto semantico all'*item* e i collocati, in questo caso, sono in larga parte nomi o verbi (Sinclair, 1993). Riprendendo l'idea di Firth, Sinclair suggerisce un'ultima interessante osservazione in un lavoro con Renouf (Sinclair & Renouf, 1991): i due studiosi pongono l'attenzione su *pattern* collocazionali in grado di dare indicazioni semantiche sui propri componenti. L'esempio più famoso è l'espressione *a X of*, dove X è, nella maggioranza dei casi, un sostantivo che indica quantità. In accordo con gli studi della *Construction Grammar*, non solo gli *item* lessicali, bensì anche la costruzione grammaticale è perciò parte della costruzione del senso. Rovesciando la prospettiva, è possibile, quindi, studiare i tipici *pattern* grammaticali di un dato *node* (e le relative relazioni con la semantica), analizzando le sue *colligazioni*. Una definizione precisa del fenomeno è data da Hoey:

Colligation can be defined as 'the grammatical company a word keeps and the position it prefers'; in other words, a word's colligations describe what it typically does grammatically (Hoey, 2000)

Nelson (2000) osserva che le colligazioni riguardano relazioni intercorrenti tra classi grammaticali, mentre le collocazioni interessano le parole appartenenti a tali classi.

Per concludere è doveroso ricordare il già citato lavoro di Mel'čuk (1998) a proposito di collocazioni e funzioni lessicali. L'innovazione di Mel'čuk è quella di ipotizzare

una tipologia universale di *lexical functions* in modo tale da suddividere il *continuum* lessicale in base a relazioni semantico-funzionali. Per funzione lessicale si intende una funzione (nell'accezione matematica) che associ a una data unità lessicale un insieme di espressioni sinonimiche sulla base della relazioni che condividono. Le collocazioni sono, in quest'ottica, legami lessicali con precise funzioni semantiche (*stark naked* di intensificatore, *a speck of dust* di quantità, ecc.). L'intero lessico può essere quindi suddiviso e organizzato in gruppi di *item* con funzioni di intensità, quantità, circostanzialità e così via in modo da coprire l'intera varietà di relazioni semantiche possibili.

In seguito allo sviluppo delle molteplici riflessioni sui legami sintagmatici esistenti tra le parole, si può dire che l'auspicio di Halliday riguardo alla concessione dello status di livello d'analisi al lessico non è stato disatteso. Esso è infatti divenuto un importante piano di indagine, strettamente interrelato a semantica e grammatica. A partire da Bally, la fraseologia si è costituita come un filone produttivo delle analisi in campo linguistico, e i lavori di Firth e della scuola neofirthiana (Halliday, Sinclair) hanno corroborato la consapevolezza dell'importanza dell'uso nella scelta degli accostamenti di parole. Sia a livello statistico-testuale che semantico, le collocazioni rimangono un terreno d'indagine stimolante e, nella loro intersezione con le unità polirematiche, uno dei principali argomenti d'interesse della contemporanea *corpus linguistics*.

L'approccio della linguistica computazionale

2.1 Calcolabilità, modelli di lingua, corpora

È chiaramente percepibile quanto il rapido e incessante sviluppo delle tecnologie informatiche durante tutta la seconda metà del XX secolo abbia prodotto grandi passi in avanti in numerose discipline e applicazioni, grazie ad algoritmi sempre più complessi e a tempi di calcolo sempre più brevi. Un tale sviluppo ha finito inevitabilmente col favorire un processo di espansione dell'informatica oltre i suoi confini naturali, strettamente legati alla matematica o all'ingegneria, per abbracciare e intersecare i settori più diversi, dove ha spesso assunto il ruolo di strumento operativo privilegiato nello studio e trattamento dei fenomeni. Il linguaggio delle macchine si è rivelato oltremodo utile nella gestione di grandi quantità di dati, permettendo, per la prima volta, di effettuare analisi i cui risultati sarebbero inevitabilmente sfuggiti al lavoro manuale dell'uomo a causa della mole di informazioni da elaborare.

Nel variegato insieme dei settori all'intersezione con l'informatica, la linguistica computazionale, che viene a costituirsi come disciplina solo a partire dagli anni '60¹, identifica oggi un importante campo di studio per un numero sempre crescente di ricerche sui fatti di lingua² e rappresenta uno dei sentieri più battuti nella ricerca teorica e applicativa sui fenomeni sintattici e lessicali collegati a polirematiche e collocazioni³.

Si è già accennato in precedenza a come la lingua storico-naturale non sia assimilabile ad un calcolo, né si conformi a modelli formali astratti, tipici dell'ottica informatica, in virtù delle imprevedibili spinte storico-sociali e delle necessità prati-

¹Nel 1962 nasce ufficialmente l'Association for Machine Translation and Computational Linguistics (AMTCL), in seguito divenuta Association for Computational Linguistics (ACL).

²Risulta impossibile tanto un esaustivo quanto un limitato elenco di riferimenti a lavori computazionali nei vasti ambiti di analisi linguistica. Si segnala, a titolo di esempio, l'imponente lavoro di Manning & Schütze (1999) che rappresenta una *summa* sulle tecniche e i metodi matematico-statistici applicabili alle analisi linguistiche.

³Si veda, a riguardo, Calzolari *et al.* (2002).

che dei parlanti⁴. Tuttavia è comunque possibile affermare che le lingue, almeno in parte, siano governate da regole che descrivano *pratiche linguistiche* (Chiari, 2007) intese come regolarità o tendenze generalmente preferite nella maggioranza degli usi. Chiari (*ibid.*) illustra come questa visione “debole” del concetto di regola linguistica sia alla base degli approcci statistico-probabilistici del trattamento della lingua, in opposizione ad una lettura in senso forte che vede la regola unicamente come regola grammaticale che produce enunciati corretti.

Sin dal principio tali due prospettive hanno visto, in ambito computazionale, una netta divisione tra i sostenitori dell'uno o dell'altro approccio. Da un lato la teoria matematica dell'informazione di Shannon & Weaver (1949) ha influenzato il dibattito computazionale in senso statistico ed induttivo, portando ai primi tentativi di modellizzazione del linguaggio attraverso le catene di Markov (Shannon, 1948), ovvero insiemi di transizioni tra “stati” possibili di numero finito⁵. Al contrario la visione opposta, supportata dalle nascenti teorie chomskyane, orientò larga parte degli studi computazionali in virtù della semplicità di formalizzazione delle regole di competenza linguistica in termini logico-matematici. In tal senso venne favorita la creazioni di modelli di tipo deduttivo, definiti *rule-based*, che stabilissero un insieme di condizioni necessarie e sufficienti alla produzione di enunciati grammaticalmente corretti.

Questo tipo di grammatica formale, tuttavia, non si è rivelata in grado di sostenere le sfide che le sempre maggiori applicazioni nel trattamento automatico della lingua richiedono. Le performance non soddisfacenti degli approcci grammaticali alla traduzione automatica e all'intelligenza artificiale⁶ hanno messo in luce la complessità insita nell'insieme di variabili integrative alle regole sintattiche di cui è pervaso ogni scambio comunicativo, quali il contesto, gli scopi comunicativi, la creatività “oltre la regola” dei parlanti, ecc. Il problema sostanziale dei modelli formali di grammatica si sintetizza nella loro visione dicotomica di ciò che risulta grammaticale o non grammaticale, in una prospettiva che potrebbe essere definita del *tutto o nulla*. Quest'ottica non contempla, ad esempio, la possibilità di usi rari di alcune parole o costruzioni in quanto l'attestazione di una frase in una produzione linguistica ne sottintende necessariamente la grammaticalità (a meno di errori).

Un'importante proprietà dei modelli statistico-probabilistici è invece la possi-

⁴A riguardo De Mauro afferma che «Proprio perché una lingua non è un'aritmetica o un'algebra, un calcolo in senso formale, ma un aggregato di forme e regole fluttuanti nel *temps* e nella *masse parlante* [...], i fatti di esecuzione (perfino quelli che possono apparire aberranti) incidono sulle unità del sistema e sulle regole, le modificano, le sopprimono, ne introducono nuove, rivelano un diverso grado di possesso degli elementi che l'intero repertorio linguistico può mettere a disposizione di una comunità a un momento dato» (De Mauro & Chiari, 2005, p. 17).

⁵Il *focus* di tali studi prevede la corretta determinazione di quale parola debba essere scelta da una macchina a continuazione di un enunciato corretto, sulla base di un numero fissato di parole precedenti prese in considerazione.

⁶Norvig (2012) specifica come al momento la totalità dei sistemi alla base di motori di ricerca, riconoscimento vocale, traduzione automatica, disambiguazione di senso siano basati su sistemi probabilistici, che garantiscono un più alto livello di output corretti.

bilità di considerare, su una scala continua, la probabilità d'uso di un'espressione, quanto essa identifichi una tendenza rara o comune della lingua e se possa rappresentare o meno una pratica linguistica prototipica secondo la definizione fornita in precedenza. Ancor più importante, rispetto a un sistema chiuso di regole astratte, è la capacità dell'approccio probabilistico di inglobare, per definizione, un insieme di *variabili nascoste* che modellizzano l'incertezza e la non conoscenza di fattori impliciti (imprevedibilmente legati a particolari caratteristiche pragmatiche degli enunciati, al contesto, a possibili errori) in modo da lasciare margine di adattabilità del modello ai dati osservati.

A riguardo è interessante considerare uno dei più famosi esempi in tema di grammaticalità adoperati nel paradigma generativo. In *Syntactic Structures* (Chomsky, 1957) vengono esposti i seguenti due esempi:

- (1) Colorless green ideas sleep furiously
- (2) Furiously sleep ideas green colorless.

Chomsky afferma che entrambi gli enunciati risultano privi di senso, ma ogni parlante della lingua inglese potrebbe stabilire con facilità che il primo, a differenza del secondo, risulta corretto grammaticalmente. Egli, tuttavia, rimprovera ai modelli probabilistici l'impossibilità di un tale discernimento:

The notion “grammatical” in English cannot be identified in any way with the notion “high order of statistical approximation to English”. It is fair to assume that neither sentence (1) nor (2) (nor indeed any part of these sentences) has ever occurred in an English discourse. Hence, in any statistical model for grammaticality, these sentences will be ruled out on identical grounds as equally ‘remote’ from English. (Chomsky, 1957)

Traspare, da questa affermazione, una visione “piatta” dei modelli statistico-probabilistici secondo cui, sulla base della non occorrenza, il modello assegni una probabilità nulla di comparsa ad eventi mai attestati prima⁷. In realtà modelli appena più complessi che considerino anche le categorie grammaticali e abbiano a disposizione un campione più o meno ampio di enunciati campione, sono in grado di astrarre le regolarità di combinazione delle parole ottenendo probabilità di cooccorrenza non nulle anche per combinazioni mai apparse nei testi del campione. Pereira (2002), ad esempio, dimostra come considerando anche modelli probabilistici poco elaborati, sia possibile trasporre in chiave statistica la nozione di grammaticalità e applicarla all'esempio di Chomsky. Attraverso un semplice modello di concatenazione di bigrammi allenato su un insieme di articoli di giornale, Pereira ottiene che

⁷È questo il caso del modello le cui stime derivino dalle frequenze relative dei soli eventi osservati (*maximum likelihood estimator*). Tale modello “ingenuo” è spesso una cattiva scelta per la tendenza sistematica a stimare per eccesso i parametri necessari all'adattamento del modello ai dati del campione (*overfitting*) (Pereira, 2002).

la probabilità di occorrenza dell'esempio (1) risulta 20.000 volte maggiore di quella dell'esempio (2). Norvig (2012) aggiunge che i modelli statistici sono in grado di attestare che entrambi gli esempi (1) e (2) risultano in ogni caso estremamente poco probabili rispetto a un enunciato "standard" come *Effective green products sell well*.

Indipendentemente dalle performance sul campo, tuttavia, il dibattito filosofico sull'accettazione dell'uno o dell'altro approccio ha radici più profonde e si ricollega alla contrapposizione di ciò che Breiman (2001) e Norvig (2012) chiamano le "due culture"⁸. Nel campo delle scienze sono venute infatti a crearsi oggi due visioni distinte che influenzano le premesse teoriche di qualsiasi indagine. La prima, definita la "cultura del *data modeling*" presuppone che la realtà e i processi che la governano siano una scatola nera che, dato un input, produce un output. La natura di tale scatola è, di base, non complessa e schematizzabile con un modello semplice, i cui parametri vanno stimati sulla base dei dati a disposizione (visione classica). Accanto a questa interpretazione esiste la "cultura del modello algoritmico", secondo cui non sempre è possibile modellizzare in maniera semplice la scatola nera della realtà. È, invece, possibile e motivato (se se ne ha la possibilità) ricorrere a complessi algoritmi e a modelli con un numero molto grande di parametri per delineare una funzione che non abbia alcuna pretesa di rappresentare la vera natura nascosta della scatola, ma il cui obiettivo sia mappare correttamente ed efficientemente i dati di input nell'output. In quest'ultimo caso, se il modello si mostra altamente predittivo, nel senso che il suo output ricalca in misura preminente le osservazioni reali a partire da dati iniziali forniti e risulta vincente nelle sfide applicative, non vi è motivo per rigettarlo, quand'anche non ricalchi ciò che Chomsky chiama l'*insight* del sistema (nel qual caso, il meccanismo del linguaggio). Chiari (2007, p. 6) fornisce la sintesi seguente:

Non è quindi necessario che il modello che usiamo per «far fare» alla macchina una determinata cosa in determinate circostanze sia lo stesso che spinge noi a un comportamento simile nell'interazione con altri esseri umani. In altri termini, è necessario un modello linguistico, ma questo modello può discostarsi dal modello linguistico usato dagli esseri umani. Allo stesso tempo è anche possibile che l'evoluzione dei modelli linguistici computazionali getti luce e fornisca suggerimenti per modelli della produzione e ricezione linguistica nelle comunità umane.

L'approccio statistico, tuttavia, non esaurisce le sue potenzialità nella sola modellizzazione dei meccanismi linguistici. L'ulteriore e fondamentale aspetto insito nella natura stessa degli studi statistici di qualsiasi tipo è l'aspirazione a *descrivere* i fenomeni sulla base dei dati empirici, individuandone le regolarità. In campo linguistico un tale approccio quantitativo fornisce inevitabilmente maggiore oggettività ed autorevolezza alle conclusioni che si avrebbero basandosi esclusivamente

⁸Breiman (2001) espone tale contrapposizione in merito alle differenze teoriche che guidano la statistica, ma essa, di fatto, si allarga a qualsiasi campo scientifico in un mondo in cui le possibilità di elaborazione delle macchine permettono rapidità di calcoli enormemente complessi.

sull'intuizione. La validità delle descrizioni statistiche della lingua è tuttavia confermata solo in presenza di un campione abbastanza vasto di dati, rappresentativo dei fenomeni che si sceglie di studiare. I *corpora*, grandi raccolte di testi, rivestono quindi un ruolo fondamentale in questo tipo di approccio, tanto sul versante del *training* dei modelli probabilistici (migliore è la base su cui il modello “impara” ad astrarre regolarità, migliore sarà la sua performance applicativa su nuovi dati in input), quanto per le statistiche quantitative sulle descrizioni linguistiche. Il grande successo contemporaneo degli approcci statistico-probabilistici (che a questo punto potremmo anche definire “empirici”) è infatti dovuto in gran parte allo sviluppo e alla crescita in dimensioni dei corpora disponibili.

Si può dire che negli ultimi decenni, infatti, la linguistica computazionale abbia subito un sempre più importante processo di avvicinamento e sovrapposizione con la linguistica dei corpora. La crescita esponenziale delle tecnologie di elaborazione e processamento di grandi quantità di testo e la loro vasta disponibilità da un lato, e le migliori performance dei modelli probabilistici dall'altro, hanno posto le basi per un connubio che ha permesso uno sviluppo sempre più accurato degli strumenti applicativi. Non meno rilevanti sono state, allo stesso tempo, le implicazioni teoriche sulla base di evidenze sperimentali che hanno spesso mostrato divergenze rispetto a ciò che in linguistica veniva considerato per assodato. Sinclair (1991, p. 4) sottolinea infatti il principio dell’“accettazione dell'evidenza” a scapito anche di più comodi costrutti mentali, precisando che:

The contrast exposed between the impressions of language detail noted by people, and the evidence compiled objectively from text is huge and systematic. It leads one to suppose that human intuition about language is highly specific, and not at all a good guide to what actually happens when the same people actually use the language.

In ragione di questo tipo di considerazioni, si è scelto di inquadrare il presente lavoro in un *framework* empirico, adottando, come si vedrà nel seguito, gli strumenti che i modelli statistico-probabilistici mettono a disposizione del linguista. Un'ottica di questo tipo ben si adatta al trattamento dei fenomeni multiparola, primariamente vari e difficilmente inquadrabili in schemi e regole formali rigide, perché fortemente condizionati dall'uso reale in cui vengono di volta in volta adottati. Il carattere di imprevedibilità delle caratteristiche proprie di tali fenomeni è, infatti, meglio analizzabile a livello quantitativo in una prospettiva *bottom-up*, in cui siano i dati empirici che scaturiscono dal corpus a guidare sistematizzazioni categoriali sul versante teorico. Allo stesso tempo la flessibilità del modello statistico appare un requisito essenziale al fine di tenere conto di una schematizzazione prototipica del comportamento e delle tendenze di un tale tipo di fenomeni, in un *continuum* centro-periferia che contempi le ineliminabili eccezioni alla regola.

Nonostante una tale premessa metodologica, polirematiche e collocazioni rappresentano da sempre, anche in ambito computazionale, un'ulteriore sfida al trattamento del linguaggio naturale. Come sottolineato da Sag *et al.* (2001), «multiword

expressions [are] a pain in the neck for Natural Language Processing». Le ragioni di tale difficoltà sono da ricercarsi nell'estrema ed eccessiva variabilità di proprietà lessicali e sintattiche delle espressioni che, per quanto evidenziate dai dati empirici, spesso mal sopportano le sistematizzazioni desiderate. L'individuazione stessa delle espressioni risulta problematica, quando effettuata dalla macchina: senza un insieme più o meno fissato di condizioni che possano identificarle, esse sono spesso indistinguibili dalle espressioni libere.

Nonostante difficoltà teoriche, pratiche e metodologiche, l'interesse ad inquadrare efficientemente polirematiche e collocazioni nel trattamento automatico della lingua non è venuto meno negli anni, soprattutto per il grande potenziale applicativo che l'approccio computazionale a questi fenomeni possiede. Non va dimenticato, inoltre, che il livello di analisi teorico non può che trarre beneficio da un approccio che metta in luce, sulla base della grande quantità di dati reali presenti nei corpora, il vero uso delle espressioni, le loro caratteristiche e i loro comportamenti specifici che sfuggirebbero alle riflessioni intuitive del linguista o, quand'anche queste ultime fossero corrispondenti al vero, mancherebbero dell'evidenza empirica che garantirebbe loro l'oggettività di un reale approccio scientifico al problema.

2.2 Principali interessi applicativi

Come accennato nel precedente paragrafo, il vivo interesse per il trattamento computazionale di polirematiche e collocazioni è fortemente motivato dal potenziale applicativo che questi fenomeni presentano in diversi ambiti dell'elaborazione del linguaggio naturale.

Il primo grande *focus* d'interesse è la possibilità di migliorare la **traduzione automatica**, attraverso una giusta corrispondenza interlinguistica tra parole grafiche e unità polirematiche o collocative, data l'importanza e la frequenza di tali fenomeni nella produzione linguistica. Le problematiche legate a questo tipo di elaborazioni si collegano alla frequente impossibilità di tradurre le espressioni parola per parola nella lingua target (si pensi alla corrispondenza it. *patata*, fr. *pomme de terre*). Spesso, infatti, a causa della più o meno forte opacità semantica e delle restrizioni distribuzionali, i traduttori dei costituenti polirematici sono imprevedibili (Monti *et al.*, 2011), non adatti a traduzioni letterali oltre i confini culturali e linguistici della lingua stessa (Barreiro, 2008) e a volte riconoscibili proprio a causa del loro *irregular translational behaviour* (Cap *et al.*, 2013) che si esplicita nella grande varietà di traduttori equivalenti quando si mettono a confronto traduzioni di riferimento in più lingue.

Un discorso di questo tipo, che sembrerebbe essere valido primariamente per le espressioni polirematiche, si estende anche alle collocazioni. Anche in presenza di legami preferenziali, infatti, in cui in generale è più facile che ogni singolo costituente abbia corrispondentemente un singolo traduttore nella lingua target, spesso non esiste predicibilità sull'uso dei traduttori a causa della non corrispondenza delle

reti collocazionali tra elementi lessicali nell'una e nell'altra lingua, come nel caso dell'italiano *prestare attenzione* e dell'inglese *pay attention*.

In generale la sfida del trattamento di polirematiche e collocazioni nella traduzione automatica sembra meglio gestibile grazie all'impiego di grandi *corpora paralleli*, vale a dire risorse di riferimento che includano gli stessi testi tradotti nelle lingue di partenza e di arrivo. Le frasi incluse nelle risorse risultano "allineate" nelle due lingue in modo da avere una corrispondenza biunivoca tra porzioni di testo di varia lunghezza. In questo modo software di tipo statistico possono essere allenati a riconoscere, attraverso un certo numero di variabili da prendere in considerazione, come i legami tra lessemi interni ad una lingua siano espressi dai legami nell'altra mediante l'astrazione di pattern sintattici e lessicali. Non va dimenticata, tuttavia, l'esistenza di altre architetture che in alcuni casi possono rivelarsi efficaci nella traduzione di costrutti idiomatici, quali i sistemi basati su regole degli approcci lessico-grammatici⁹, come mostrato in Monti *et al.* (2011).

Sempre in ambito traduttivo, la specificità di connessioni lessicali che ogni lingua sviluppa attraverso le collocazioni risulta, spesso, un argomento delicato anche per gli stessi traduttori umani, la cui presenza è ancora fortemente richiesta in virtù del fatto che gli attuali sistemi di traduzione automatica non sono ancora perfetti (Barrachina *et al.*, 2009). Tuttavia eventuali interferenze con la lingua L1, l'uso altamente peculiare in particolare di preposizioni, terminologie, accezioni specialistiche e costrutti altamente idiomatici possono diminuire il grado di accuratezza e di naturalità della traduzione prodotta. Un filone degli studi in linguistica computazionale si è dedicato, negli anni, allo sviluppo di tecnologie per la **traduzione assistita dal computer**, vale a dire un insieme di risorse informatiche fruibili dai traduttori, che adiuvasse il loro lavoro, il cui più noto esempio è quello delle *memorie di traduzione* (Bowker, 2002; Somers, 2003). In questo settore il trattamento computazionale di collocazioni e polirematiche (specie in ambito terminologico) risulta fondamentale, ad esempio quando si voglia rendere efficiente la corrispondenza tra i segmenti tradotti custoditi nelle memorie di traduzione e le entità lessicali in questione. L'accessibilità a risorse informatiche che, tenendo presente i legami polirematici e collocazionali, suggeriscano le scelte migliori al traduttore, permetteranno infatti di migliorare una già elevata competenza linguistica di partenza. Quest'ultima, infatti, «è caratterizzata [...] soprattutto dalla capacità di utilizzare le combinazioni lessicali *proprie* della lingua», poiché «una ridotta abilità nell'uso delle collocazioni produce un linguaggio povero, incompleto e poco articolato che, a sua volta, comprometterà la comunicazione molto più significativamente di quanto non faccia un'espressione che presenti carenze sintattiche e grammaticali» (Tiberii, 2012, p. 3).

Un discorso di questo tipo, tuttavia, si ricollega all'altrettanto fondamentale que-

⁹Benché fondati su regole, gli approcci nell'ambito della lessico-grammatica formalizzano pattern sintattico-semantici sulla base dell'analisi di grandi quantità di dati empirici, evitando in questo modo le ottiche formali *hardcore* del paradigma generativo.

stione della necessità di vaste banche dati lessicografiche che custodiscano in maniera strutturata le informazioni relative ai fenomeni polirematico-collocazionali. Sia che si necessiti di un uso specialistico da parte di traduttori, che di una consultazione generale ad opera di un utente comune, la disponibilità di risorse più o meno esauritive che annoverino tra i lemmi 'parole formate da più parole' o che considerino le preferenze lessicali di cooccorrenza dei vari lessemi, appare come una necessità primaria nell'ambito della **lessicografia** contemporanea (Heid & Weller, 2010).

Gli approcci computazionali forniscono un supporto strategico alla creazione di dizionari che includano i fenomeni multiparola. I metodi di riconoscimento ed estrazione automatica di espressioni che presentino un alto grado di coesione facilita il censimento di questo tipo di fenomeni, specie sul versante terminologico. Non sono, infatti, solo i dizionari (mono- e bilingui) del lessico generale a beneficiare dell'inclusione delle polirematiche tra i lemmi dell'opera, bensì primariamente i dizionari dei lessici tecnico-specialistici, in cui le polirematiche formano tipicamente terminologia. Nonostante molti settori presentino ormai un lessico terminologico ampiamente catalogato, i settori di ricerca in rapida evoluzione, il più delle volte a carattere multidisciplinare, tendono a sviluppare termini nuovi e vari, non inclusi in liste di riferimento o altre fonti autorevoli (Heid & Gojun, 2012). Diventa di estremo rilievo, a questo punto, lo sviluppo di strategie atte a identificare e lemmatizzare questo tipo di espressioni, anche e soprattutto per la traduzione (automatica, assistita o manuale), specie in mancanza di grandi testi di dominio tradotti e presenti in corpora paralleli.

Parallelamente, gli sforzi computazionali sono mirati a sviluppare anche basi dati di collocazioni, quali fenomeni di cui rendere conto a livello lessicografico. Come ricorda Cop (1991, p. 2776), le collocazioni, anche se trasparenti a livello semantico, necessitano l'inclusione nei dizionari semplicemente perché «they are not predictable». In questo senso sono venuti alla luce, negli anni, diversi dizionari combinatori sia in ambito anglosassone¹⁰ che in ambito romanzo, con contributi ad esempio sullo spagnolo¹¹ e l'italiano¹² che mirano ad aiutare sia i traduttori che gli utenti comuni nella scelta più adatta delle combinazioni lessicali proprie della lingua.

2.3 Collocazioni empiriche

Fin da subito gli studi su contesto e cooccorrenza in linguistica computazionale evidenziano come esista una particolare attrazione tra determinati lessemi in una lingua e per la prima volta tale specificità può essere quantificata in termini di frequenza di cooccorrenza in un corpus¹³. La principale proprietà, infatti, che

¹⁰Si veda, ad esempio, il lavoro di McIntosh *et al.* (2009), rivolto principalmente a studenti dell'inglese come L2.

¹¹Si veda Bosque (2004b).

¹²Si vedano i lavori di Urzì (2009), Lo Cascio (2011), Tiberii (2012).

¹³Cfr. par. 2.4.1.

espressioni come *heavy smoker* o *black box* esibiscono è una maggiore frequenza di occorrenza congiunta dei componenti, rispetto alla frequenza con cui questi ultimi si presentano nel contesto di altri lessemi.

Nella prospettiva firthiana, tuttavia, l'alta ricorsività di determinati gruppi di parole rimane unicamente sul piano empirico, non presupponendo implicazioni di carattere fraseologico. È questa un'importante differenza che, come sottolineato in Evert (2008), ha spesso generato confusione rispetto alla nozione di collocazione nel senso neofirthiano di combinazione lessicalmente determinata o a quella relativa alle combinazioni lessicalizzate di parole che mostrano specificità sintattiche e semantiche che in questo lavoro vengono chiamate polirematiche.

Nonostante il significato e l'uso di una parola vengano influenzati e determinati dai componenti del contesto, la definizione empirica di collocazione quale frequente cooccorrenza di combinazioni di parole lascia ampio spazio ad associazioni come ad es. *notte/giorno*, *gatto/miagolare* che, pur presentandosi spesso vicine, non sembrano identificare legami lessicali o specificità semantiche dell'espressione se non, a volte, nel senso coseriano di "solidarietà" (cfr. par. 1.5.4). Anche Bosque (2004a, p. LXXXIV) sottolinea le problematicità legate al concetto di frequenza di occorrenza in ambito collocazionale:

En realidad, la frecuencia no es un factor que garantiza la idiomatidad, y tampoco la rutina o el cliché en ninguna de sus formas posibles. La combinación del verbo *leer* y el sustantivo *libro* es sumamente frecuente. Si no nos fijamos de nuestra introspección, podemos buscarla en los corpus, contar las veces que aparece y comparar esa proporción con las que corresponden a otros sustantivos que se combinan con el verbo *leer* y a otros verbos que se combinan con el sustantivo *libro*. Supongamos que hacemos todo eso y que observamos que esas frecuencias son altas. ¿Qué hemos descubierto? Me parece que la respuesta correcta es "Prácticamente nada".

È doveroso precisare che tali associazioni, pur rappresentando dei fenomeni che, secondo alcuni, non necessitano di uno status categoriale nella teoria linguistica, esprimono una forte capacità predittiva su quali lessemi possono accompagnare quali altri e non è escluso che nel tempo la ricorrente vicinanza di una parola all'altra tenda a specializzare o creare nuove accezioni tra i componenti e a fondare in tal senso una collocazione lessicale.

Al fine di evitare ulteriore confusione in un già delicato ambito definitorio, si definirà *collocazione empirica* (in senso firthiano) una semplice combinazione di parole particolarmente ricorrente o statisticamente rilevante, la cui frequenza è attestata in un corpus.

È evidente che, in questo senso, collocazioni empiriche da un lato e collocazioni (lessicali) e polirematiche dall'altro formano due insiemi di entità intersecantisi ma non sovrapponibili, come schematizzato in Figura 2.1.

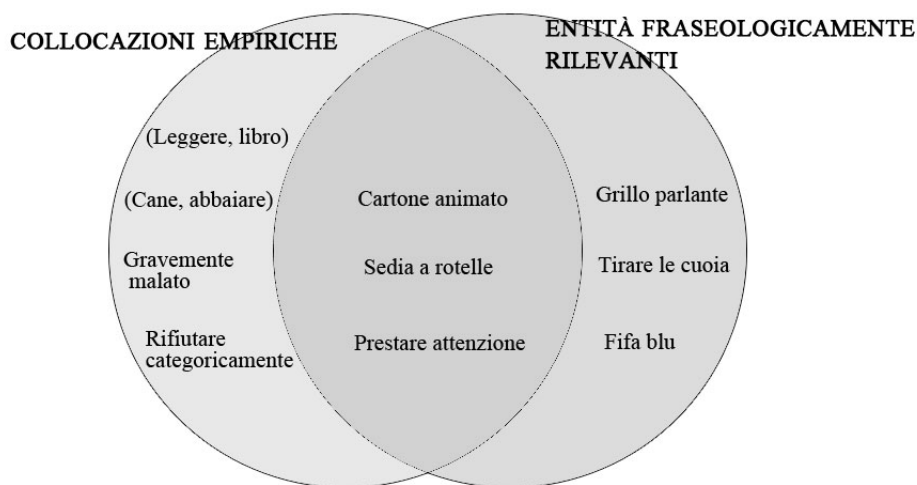


Figura 2.1: Schema del rapporto tra collocazioni empiriche ed entità fraseologicamente rilevanti, che includono unità polirematiche e collocazioni. Dal corpus PAISÀ (Lyding *et al.*, 2014), composto di 250 milioni di parole, si ottengono le seguenti frequenze per le espressioni riportate: (Leggere, libro) 1934; (Cane, abbaiare) 90; Gravemente malato 288; Rifiutare categoricamente 124; Cartone animato 3038; Sedia a rotelle 746; Prestare attenzione 503; Grillo parlante 54; Tirare le cuoia 24; Fifa blu 3. I dati si riferiscono a query su componenti lemmatizzati; per le coppie di componenti tra parentesi la ricerca è stata effettuata su tutte le flessioni dei componenti includendo eventuali parole intervenienti. Una soglia di frequenza pari a 60 esclude dalle collocazioni empiriche le espressioni più a destra.

È ragionevole, infatti, ritenere che molte (ma non tutte) delle combinazioni che con buona probabilità ricorrono con frequenza alta in un corpus siano lessemi complessi e, allo stesso tempo, è possibile che molte collocazioni lessicali o polirematiche non compaiano così frequentemente nell'uso da apparire come espressioni ricorrenti.

In relazione alla prima ipotesi, infatti, Evert (2008, p. 1214) chiarisce molto bene che:

There is a close connection between empirical collocations and multi-word expressions, though. A thorough analysis of the collocations found in a corpus study will invariably bring up non-compositionality and lexicalisation phenomena as an explanation for many of the observed collocations.

In questo modo la collocatività empirica viene individuata come un epifenomeno derivante da fattori lessicali e semantici.

Riguardo al secondo punto ancora Evert (*ibid.*) sottolinea come uno degli idioms più noti e maggiormente preso ad esempio negli studi anglosassoni, *kick the bucket*, sia in realtà un'espressione che in una risorsa come il British National Corpus (che conta circa 100 milioni di parole) compare solo 3 volte. Esempi di questo tipo mettono in luce un aspetto cruciale negli studi linguistici su collocazioni e polire-

matiche, e cioè che la frequenza, da sola, non è sufficiente a mettere in risalto la collocatività di un'espressione. Come si vedrà nel seguito, per ovviare a questo problema l'"attrazione" tra i componenti di un'espressione viene oggi di solito espressa mediante una misura d'associazione.

Va infine sottolineato come in generale, a livello metodologico, ogni studio empirico non possa che focalizzarsi sul solo insieme delle collocazioni empiriche. Ogni fenomeno multiparola che non vi rientri ha come causa dell'assenza solo due possibilità: esso può non essere presente nel corpus considerato, oppure presentare un basso numero di occorrenze, il che lo rende non individuabile attraverso un filtro di frequenza o attraverso le misure d'associazione che spesso dipendono fortemente dalla frequenza stessa. Quand'anche espressioni di bassissima frequenza venissero prese in considerazione nello studio di comportamenti generali relativi ai fenomeni multiparola, sorgerebbe allora il rischio di avallare conclusioni poco attendibili dovute alla bassa rappresentatività che il numero esiguo di occorrenze (e quindi di comportamenti, variazioni possibili e contesti di occorrenza) delle espressioni considerate può garantire. Uno studio rivolto ad un esame del fenomeno nella sua generalità, quindi, nel caso di un campione di collocazioni empiriche selezionato attraverso una metodologia attendibile e ragionata, dovrà assumere che quest'ultimo sia rappresentativo e che le conclusioni ad esso relative siano espandibili anche alle entità fraseologiche non prese in considerazione nell'analisi¹⁴.

2.4 Il trattamento informatico e statistico delle parole

2.4.1 Definizioni e concetti generali

L'idea di compiere analisi e studiare la lingua e i suoi fenomeni grazie a dati quantitativi che possano essere calcolati ed estratti da corpora di vario tipo presuppone che i testi che costituiscono la base empirica della ricerca siano presenti in una qualche memoria informatica sotto forma di *file* interrogabili. Ciò è possibile in quanto, potenzialmente «ogni informazione linguistica può essere catturata in maniera automatica, essendo costituita da una sequenza di *bytes* (ovvero una catena di caratteri)» (Bolasco, 1999, p. 180), si tratti delle parole di un testo scritto, della trasposizione di un testo orale o di metainformazioni riguardo contesto, scopi comunicativi, atti paralinguistici, prossemica, ecc. La trasposizione di un testo in formato elettronico, per quanto all'apparenza possa sembrare un processo intuitivamente semplice, nasconde una lunga serie di operazioni di codifica (di norma convenzionalizzate) che hanno lo scopo di collegare l'idea di *parola* visualizzabile sullo schermo ai bit che ne costituiscono la forma base in linguaggio macchina.

¹⁴A riguardo si rimanda alle considerazioni del par 3.3.1.

Di tutti i caratteri identificabili elettronicamente e disponibili nella composizione di un testo, per convenzione ve ne è una parte che non viene assimilata all'alfabeto, definita come l'insieme dei *separatori*, di cui alcuni esempi possono essere lo spazio bianco, il tabulatore, i segni di punteggiatura, le parentesi, ecc. Per contro, ciò che tra i possibili caratteri non viene incluso tra i separatori diviene un *carattere alfabetico*. In ambito computazionale, per convenzione, una parola viene quindi definita come qualsiasi sequenza di caratteri alfabetici delimitata agli estremi da due separatori¹⁵. Più specificamente, questa definizione di parola esprime ciò che viene generalmente indicato con il termine *occorrenza* (*token* in inglese). Un'ulteriore definizione interessa tutte le occorrenze composte dalla stessa sequenza di caratteri, che possono essere identificate come appartenenti allo stesso *type*¹⁶. Infine è possibile astrarre una categoria più ampia, detta *lemma*, che racchiude tutti i *type* (e di conseguenza le relative occorrenze) riconducibili alla forma di citazione convenzionale della parola nei dizionari (che quindi comprende tutte le flessioni)¹⁷. Le diverse varianti che costituiscono i *type* riconducibili ad uno stesso lemma vengono dette *forme*. Il numero di occorrenze raggruppate sotto uno stesso *type*, lemma o forma, rappresenta la *frequenza* della categoria considerata.

È importante sottolineare che esistono strumenti computazionali in grado di assegnare in maniera automatica il lemma di riferimento a tutti i token di un corpus, mediante l'uso di dizionari, alberi decisionali e della statistica delle probabilità. In generale questo processo viene accompagnato anche dall'assegnazione delle categorie grammaticali ai relativi token, come nel caso del noto *tool* TreeTagger (Schmid, 1994), e si parla in questo caso di un processo di *part-of-speech tagging* (*POS-tagging* nel seguito).

In relazione ai fenomeni multiparola, assume un'importanza rilevante la definizione di *cooccorrenza*, ovvero l'attestazione contemporanea e vicina, in un testo, di due o più componenti di un'espressione. La definizione di cooccorrenza risulta fondamentale in relazione alla metodologia applicata nello sviluppo delle analisi computazionali su polirematiche e collocazioni in quanto da essa dipende la diversa quantificazione dei dati di frequenza di espressione e componenti (Evert, 2008), che a loro volta costituiscono la base per la misura di associazione lessicale tra i componenti.

In prima istanza è possibile definire un tipo di cooccorrenza cosiddetta '*superficiale*', nell'ipotesi in cui il corpus a disposizione non abbia subito un processo di *parsing* sintattico, ovvero non siano stati esplicitati i legami di dipendenza sintattica (eventualmente anche a distanza) tra i componenti della frase. In questo caso

¹⁵Cfr. la definizione di *parola grafica* vista al par. 1.1.

¹⁶Nella frase 'I bravi professori formano dei bravi studenti' il *type* 'bravi' conta due occorrenze. Va sottolineato che anche le variazioni di lettere maiuscole o minuscole contribuiscono a formare *type* distinti, in quanto rappresentati da sequenze diverse di caratteri, come nei casi di *bello* e *Bello*.

¹⁷Casi del tipo *amò, amai, amato, amassi* o *bei, bel, belle* sono diversi *type* riconducibili, rispettivamente, ai lemmi *amare* e *bello*.

l'unica possibilità di analisi della cooccorrenza si ricollega alla vicinanza reciproca delle forme occorrenti.

Il caso più semplice è quello per cui una cooccorrenza diventa la sequenza di due o più forme contigue e, detto n il numero di componenti considerati per un'espressione, si parla in questo caso di n -grammi (*bigrammi* nel caso di due componenti, *trigrammi* per tre componenti, e così via). In Figura 2.2 sono illustrate alcune occorrenze del bigramma *cartone animato* tratte dal corpus PAISÀ (Lyding *et al.*, 2014).

- (1) La sua fama si consolidò negli anni Trenta con una serie di opere fra cui la famosa serie di Capitano Futuro, da cui fu poi realizzato un cartone animato giapponese.
- (2) Joe Casey è anche sceneggiatore televisivo (suo è il cartone animato Ben 10, trasmesso da Cartoon Network) ed è il leader del gruppo indy rock Sellouts.
- (3) Nel cartone animato I Griffin interpreta sé stesso nella carica di sindaco della città immaginaria di Quahog.

Figura 2.2: Esempi di occorrenza del bigramma *cartone animato* (sottolineato nel testo) nel corpus PAISÀ.

La fortuna delle analisi basate su n -grammi è dovuta in larga parte al loro utilizzo nei modelli probabilistici di generazione di enunciati attraverso catene di Markov, secondo cui la probabilità di utilizzo e attestazione di una certa parola è strettamente legata alle parole che immediatamente la precedono¹⁸. Non mancano, tuttavia, diversi lavori in letteratura che hanno preso in considerazione cooccorrenze di questo tipo per una analisi computazionale sui fenomeni multiparola in virtù della semplicità della loro formalizzazione operativa, che non presuppone alcuna elaborazione sul testo in termini di ricerca di un qualche legame tra i componenti che vada oltre la mera contiguità¹⁹.

Più in generale, la “vicinanza” delle forme cooccorrenti può oltrepassare la semplice contiguità dell' n -gramma, per comprendere i casi in cui le forme che generano la cooccorrenza si trovano separate entro un certo numero di parole intervenienti. La distanza che separa le forme cooccorrenti, in termini di quantità di parole intervenienti, viene denominata *span*, ed essa può essere arbitrariamente scelta a seconda dell'analisi che si intende portare a termine, nonostante, di norma, non oltrepassi i confini della frase. Un'indagine su fenomeni collocativi privilegerà, in genere, *span* variabili da 3 a 5 parole (Sinclair, 1991) come nell'esempio di Figura 2.3, mentre intervalli di parole di gran lunga maggiori, che possono espandersi al limite della frase stessa (come negli esempi di Figura 2.4) possono rivelarsi utili in studi di semantica distribuzionale (Schütze, 1998) o nell'analisi delle corrispondenze lessicali (Bolasco, 1999, 2013).

¹⁸Per una panoramica generale si veda il lavoro di Damerau (1971).

¹⁹Come ricordato in Evert (2008), tecniche di approccio alla collocatività basata su bigrammi risalgono a Choueka (1988) o anche a Schone & Jurafsky (2001), mentre non meno importante risulta il più recente lavoro di Ramisch *et al.* (2010b) sull'identificazione automatica basata su misure d'associazione statistica applicate a n -grammi.

- (1) Il partito governante Russia Unità ha un ampio supporto popolare e i media controllati dal governo hanno prestato poca attenzione alle opposizioni al partito.
- (2) Occorre prestare la massima attenzione al morso di un cane a causa della possibilità di contrarre una grave malattia: la rabbia o idrofobia.
- (3) È opportuno prestare attenzione ai bordi della lastra, che terminano a strapiombo ai limiti del bosco e della prateria, senza alcuna protezione.

Figura 2.3: Esempi di occorrenza della collocazione *prestare attenzione* con uno span massimo di 3 parole nel corpus PAISÀ.

- (1) Nel 1945, mentre la seconda **guerra** mondiale stava per concludersi, la Germania Nazista venne “incastrata” tra gli eserciti degli alleati occidentali, che dall’ovest appunto giungevano, e l’Unione Sovietica che con i suoi **eserciti** avanzava da est.
- (2) In quegli anni era in corso la **guerra** di successione spagnola alla quale poneva fine la vittoria dell’**esercito** sabaudo sulla Francia nella battaglia di Torino del 6 settembre 1706.
- (3) Cleopatra, però, radunò un **esercito** ed iniziò una **guerra** civile in Egitto.

Figura 2.4: Esempi di cooccorrenze della coppia di lemmi (*guerra, esercito*), riportati in grassetto, all’interno di frasi del corpus PAISÀ. La compresenza di lemmi in frammenti di vario tipo, tra cui le frasi, viene spesso sfruttata al fine di creare modelli semantici della distribuzione del lessico.

Una seconda categoria di cooccorrenze viene invece definita non per la vicinanza testuale, bensì in relazione al legame sintattico che lega i componenti, e per questo esse prendono il nome di *cooccorrenze sintattiche*. Al fine dell’identificazione dei componenti come parte di una cooccorrenza essi devono quindi esibire un legame diretto in termini di dipendenza²⁰: informazione, questa, generalmente disponibile se il corpus in analisi ha subito un processo di *parsing* sintattico, che ha quindi riconosciuto mediante indici puntatori alle varie forme i modificatori, le entità modificate, le reggenze, ecc. Gli aggettivi sono quindi ricondotti ai sostantivi che modificano, soggetto ed oggetto sono ricondotti al verbo da cui dipendono, e così via. Come ricordato in Evert (2008), le cooccorrenze sintattiche risultano particolarmente indicate nel trattamento di lingue che esibiscono dipendenze sintattiche tra componenti a grande distanza. In Figura 2.5 sono riportati alcuni esempi di cooccorrenze sintattiche.

²⁰Ad oggi, infatti, l’annotazione sintattica preferita nel trattamento computazionale è il *dependency parsing*. La nozione di dipendenza presuppone che tutte le parole in una frase siano connesse tramite collegamenti diretti. Il verbo costituisce generalmente il centro strutturale della frase, mentre ogni altra parola dipende da esso in maniera diretta o indiretta. La struttura delle dipendenze, inoltre, differisce dalla grammatica chomskyana dei sintagmi poiché essa non prevede nodi sintattici. La letteratura sul tema risulta sconfinata per poter essere approfondita in questa sede, tuttavia si rimanda al lavoro di base di Tensiè (1959) per un’impostazione teorica, nonché all’ottima sintesi di Nivre (2005) per un percorso storico fino ad arrivare agli approcci computazionali moderni.

-
- (1) C'è anche la possibilità di mettere alla prova il proprio **pollice verde**, cimentandosi nelle attività di semina e trapianto delle erbe officinali.
- (2) Le temperature sono previste in **lieve e generale rialzo** su tutte le regioni.
- (3) Migliaia di start-ups nei nuovi settori dell' high-tech-dot-com hanno **tirato le cuoia**

Figura 2.5: Esempi di cooccorrenze sintattiche, riportate in grassetto, all'interno di frasi del corpus PAISÀ.

2.4.2 Frequenze di cooccorrenza

Qualsiasi sia la tipologia di cooccorrenza presa in esame nello studio del lessico di un corpus, al fine di individuare tra le innumerevoli combinazioni di parole quali siano i componenti legati da un qualche legame particolare (ovvero che costituiscano collocazioni empiriche), è necessario scegliere un indice che quantifichi la 'collocatività' fra parole. In altri termini è necessaria una qualche misura che, sulla base di informazioni estraibili dal corpus, sia in grado di selezionare un insieme più o meno ristretto di espressioni d'interesse o, quanto meno, riesca a fissare due poli di un continuum che vedano ad un estremo le espressioni più tipicamente fisse e dall'altro le espressioni generalmente libere.

La preponderanza di due o più parole ad occorrere in maniera congiunta può essere indice del fatto che esse siano parte di un'espressione 'cristallizzata' o meno libera rispetto a combinazioni di lessemi che possono essere accostati con facilità ad un numero elevato di altre parole casuali. In tal caso si potrebbe, quindi, considerare la frequenza di cooccorrenza come una misura atta ad indentificare, nell'insieme delle innumerevoli cooccorrenze del corpus, quelle che con maggiore probabilità ricadono nei fenomeni multiparola. Tuttavia è facilmente verificabile che una lista di espressioni ordinata in maniera decrescente per frequenza di cooccorrenza, pur includendo espressioni che appartengono alla sfera delle unità polirematiche o delle collocazioni, selezionerà un gran numero di espressioni generalmente considerate libere, come mostrato in Tabella 2.1. Si può dire, in altri termini, che la frequenza di cooccorrenza ha una bassa *precisione* in molti casi.

La ragione dell'insufficienza di tale metodologia è dovuta al fatto che, in generale, non è possibile stabilire se l'alto numero di cooccorrenze dei componenti sia il risultato di una loro particolare attrazione o del fatto che essi, essendo già molto frequenti di per sé, hanno grande possibilità di combinarsi insieme. Un esempio magistrale di quest'ultimo caso è riportato, ad esempio, in Evert (2008), dove si prende in esame il bigramma inglese *is to*, analizzandone le occorrenze nel Brown corpus (Francis & Kucera, 1964). In tale risorsa, composta da circa un milione di parole, i componenti *is to* presentano 260 cooccorrenze, il che identifica il bigramma come uno dei più frequenti del corpus. Tuttavia gli stessi componenti, considerati singolarmente, esibiscono un alto numero di occorrenze: *is* ha una frequenza di occorrenza

Bigrammi		N + A		V + Adv	
espr.	freq.	espr.	freq.	espr.	freq.
per il	673.789	collegamento esterno	131.642	essere sempre	17.446
e il	639.965	guerra mondiale	27.803	essere ancora	13.658
con il	609.485	anno successivo	18.341	essere più	12.477
il suo	479.968	personalità legata	17.233	essere quindi	8.543
essere il	419.305	colonna sonora	9.197	essere invece	8.322
essere essere	495.467	serie televisiva	9.064	essere infatti	7.577
che il	326.542	anno seguente	7.296	essere così	7.340
e di	297.532	chiesa cattolica	7.127	essere già	7.028
tutto il	293.978	centro storico	6.787	essere inoltre	6.886
di un	276.850	essere umano	6.704	fare sì	6.792

Tabella 2.1: Prime dieci espressioni più frequenti nel corpus PAISÀ appartenenti a diversi pattern. Nella prima tabella sono inclusi i bigrammi (dati di cooccorrenza superficiale), nella seconda e nella terza le combinazioni di nome e aggettivo e di verbo ed avverbio (dati di cooccorrenza sintattica). In tutti i casi le query considerano i componenti lemmatizzati.

pari a circa 10.000, mentre *to* viene attestato circa 26.000 volte. Nell'ipotesi in cui le due parole si combinassero in maniera del tutto casuale, si avrebbe che per ognuna delle 10.000 occorrenze di *is*, ci sarebbe una probabilità pari a 26.000 su un milione (cioè 26 su 1.000) che la parola successiva sia *to*. In questo caso, il numero di volte per cui il bigramma dovrebbe comparire nel corpus dovrebbe risultare pari a $10.000 \cdot (26/1000)$, ovvero 260. La cosiddetta frequenza di cooccorrenza *aspettata* nel caso di combinazione casuale, risulta uguale alla frequenza di cooccorrenza *osservata* nel corpus. In questo caso, quindi, l'alta frequenza di *is to* è frutto di una pura casualità e non indice di un legame collocativo: non esiste, cioè, un'*associazione* tra i componenti. Considerazioni di questo tipo suggeriscono che una misura quantitativa che sia in grado di evidenziare legami tra i componenti di un'espressione debba non solo prendere in considerazione la frequenza di cooccorrenza, ma anche la frequenza di occorrenza dei singoli componenti, che viene definita *frequenza marginale*. Grazie alle frequenze marginali, infatti, si è in grado di risalire alla frequenza attesa in caso di combinazione casuale attraverso una formula standard. Nel caso di due soli componenti, ad esempio, detta f_1 la frequenza del primo componente, f_2 la frequenza del secondo ed N il numero totale di parole del corpus, la frequenza attesa dell'espressione risulta:

$$F_{aspettata} = \frac{f_1 f_2}{N} \quad (2.1)$$

La frequenza attesa di cooccorrenza formalizza, quindi, la situazione di *indipendenza statistica* tra i componenti, anche detta *ipotesi nulla*, rappresentando il numero di cooccorrenze dei componenti che attesta la loro non associazione.

Avendo a disposizione, come dato empirico estraibile dal corpus, la frequenza osservata dell'espressione, e potendo calcolare la frequenza attesa, è possibile costruire delle *misure d'associazione*, ovvero degli indici che, in diverso modo, mettono a confronto le frequenze suddette o quantità che dipendono da esse. Maggiore risul-

terà la discrepanza tra le due, maggiore sarà la probabilità che l'espressione includa parole che non si combinano insieme per puro caso e che per questo possa essere riconducibile ai fenomeni multiparola.

2.4.3 Misure d'associazione

Per quanto detto nel precedente paragrafo, una misura d'associazione si configura come una formula che, in base ai dati statistici ricavati dal corpus, è in grado di calcolare un punteggio d'associazione fra componenti per ogni espressione considerata, ovvero un numero che possa essere utilizzato (a) per ordinare in base alla significatività di associazione le espressioni in maniera decrescente (più è alto il punteggio di associazione, più l'espressione ha probabilità di rientrare nei fenomeni multiparola, cfr. Tabella 2.2) o (b) per selezionare un insieme finito di espressioni fissando una soglia limite (tutte le espressioni con un punteggio più alto della soglia sono collocazioni empiriche).

La letteratura sulle misure d'associazione atte a individuare in maniera automatica le collocazioni empiriche si è rivelata fiorente, specie negli ultimi anni. Tuttavia l'idea di applicare i metodi statistici al linguaggi naturale per ottenere una quantificazione dell'attrazione fra parole era nata già negli anni '60, in concomitanza con la nascita della linguistica computazionale: in Stevens *et al.* (1965), ad esempio, è possibile già ritrovare più di dieci differenti misure d'associazione, sebbene ancora senza fondati criteri d'interpretazione e valutazione dei risultati (Giuliano, 1965, p. 259).

espressione	punteggio
medaglia d'oro	0,0331
corso d'acqua	0,0184
luogo d'interesse	0,0174
opera d'arte	0,0169
medaglia d'argento	0,0146
lunghezza d'onda	0,0132
storia d'amore	0,0116
diritto d'autore	0,0115
televisione via cavo	0,0111
tempesta d'amore	0,0106

Tabella 2.2: Le 10 espressioni più rilevanti per punteggio d'associazione statistica del pattern NPN nel corpus PAISÀ. La misura d'associazione considerata nell'esempio è il *coefficiente di Dice* (Smadja *et al.*, 1996; Dias *et al.*, 1999), la cui espressione matematica è data in Figura 2.7. I dati sono stati ottenuti su trigrammi conformi al pattern considerato grazie al software *mwetoolkit* (Ramisch *et al.*, 2010b). La query considera componenti lemmatizzati ed espressioni con almeno 69 occorrenze nel corpus.

Il problema della mancanza di grandi quantità di testi disponibili in formato elettronico da un lato, e dei tempi di calcolo degli algoritmi su grandi moli di dati dall'altro, non ha permesso, se non a partire dagli anni '90, lo sviluppo di una vera e propria letteratura sul tema. Compaiono solo in questo decennio, infatti, alcune tra le misure d'associazione di maggior successo nelle analisi linguistiche e lessicografiche: *mutual information* (Church & Hanks, 1990), *t-score* (Church *et al.*, 1991) e *log-likelihood* (Dunning, 1993) (le cui formule di calcolo sono espresse in Figura 2.6) insieme ad un vasto numero di altri indici più o meno noti.

$$MI = \log_2 \frac{O}{E} \quad (2.2)$$

$$\text{simple-log-likelihood} = 2 \left[O \log \frac{O}{E} - (O - E) \right] \quad (2.3)$$

$$\text{t-score} = \frac{O - E}{\sqrt{O}} \quad (2.4)$$

Figura 2.6: Tre delle più usate misure d'associazione. Nelle formule O rappresenta la frequenza osservata dell'espressione, E la frequenza attesa.

Durante gli stessi anni, inoltre, la valutazione della bontà delle performance delle misure d'associazione varia molto a seconda delle pubblicazioni e degli obiettivi di ricerca. I lavori di matrice più prettamente linguistica tendono a focalizzarsi su casi di studio limitati tentando, come sottolineato in Evert (2004, p.29), di ottenere una comprensione intuitiva delle differenze tra le misure, sulla base dell'evidenza dei dati di estrazione²¹ o focalizzandosi su descrizioni delle caratteristiche di singole misure²². Lavori, invece, più orientati alla statistica riescono a compiere valutazioni maggiormente obiettive delle performance di alcune misure d'associazione su vasta scala, riuscendo a fornire motivazioni sulle criticità di alcune misure in determinate situazioni²³.

Tuttavia, solo nel decennio successivo i lavori di Evert (2004, 2008) riusciranno a configurarsi come una *summa* delle misure d'associazione statistica disponibili, analizzate in una prospettiva rigorosa sia dal punto di vista matematico che dal punto di vista delle performance.

Diventa chiara la differenza tra misure d'associazione *semplici*, che prendono in considerazione unicamente frequenza attesa e osservata dell'espressione, e misure

²¹In Lapata *et al.* (1999), ad esempio, si considerano i punteggi dati dalle misure d'associazione a diverse espressioni con i giudizi di plausibilità dell'associazione forniti da parlanti madrelingua. Stubbs (1995) si sofferma su cooccorrenze statisticamente significative di un ristretto numero di lemmi specifici. Il famoso compendio di Manning & Schütze (1999) si limita a discutere delle misure d'associazione su liste di pochi «interesting bigrams» (Evert, 2004, p. 30).

²²Cfr. Smadja *et al.* (1996).

²³Si vedano, ad esempio, Krenn (2000), Evert & Krenn (2001).

statistiche (cfr. la Figura 2.7), che utilizzano, invece, *tabelle di contingenza* e che risultano più rigorose.

$$\text{average-MI} = \sum_{ij} O_{ij} \log_2 \frac{O_{ij}}{E_{ij}} \quad (2.5)$$

$$\text{log-likelihood} = 2 \sum_{ij} O_{ij} \log \frac{O_{ij}}{E_{ij}} \quad (2.6)$$

$$\text{Dice} = \frac{2O_{11}}{R_1 + C_1} \quad (2.7)$$

Figura 2.7: Alcune delle misure d'associazione statistica più utilizzate. Log-likelihood e average-MI risultano gli equivalenti statistici delle misure semplici di Figura 2.6. T-score non è presente in quanto il calcolo dell'associazione semplice si basa su una procedura empirica non giustificata matematicamente e non può, quindi, avere un'equivalente statistico (Evert, 2008, p. 1235). Il coefficiente di Dice, non presente tra le misure semplici, ha un significato unicamente statistico.

Le tabelle suddette riassumono, per ogni coppia di componenti²⁴, i dati relativi alle frequenze aspettate e osservate, oltre che dell'espressione, anche delle combinazioni dei singoli componenti con tutte le parole del corpus e delle ulteriori combinazioni in cui nessuno dei componenti è presente, come mostrato in Figura 2.8.

A seconda dell'interesse della ricerca, che necessita una maggiore o minore rigorosa trattazione statistica, è possibile optare per l'una o l'altra tipologia di misura.

Negli stessi lavori, inoltre, le differenti metodologie adottate dalle diverse misure vengono ricondotte a due grandi gruppi sulla base di criteri presi in considerazione nel computo del punteggio d'associazione, e cioè misure di *effect size* e misure di *significatività*:

The former [measures of *effect size*] ask the question “how strongly are the words attracted to each other?” (operationalised as “how much does observed cooccurrence frequency exceed expected frequency?”), while the latter [measures of *significance*] ask “how much evidence is there for a positive association between the words, no matter how small effect size is?” (“operationalised as “how unlikely is the null hypothesis that the words are independent?”) (Evert, 2008, p. 1228).

In altri termini, le misure di *effect size* (tra le quali, ad esempio, la *mutual information*) forniscono una stima diretta dell'associazione attraverso la magnitudine della discrepanza tra frequenze osservate e aspettate. Le misure di significatività (tra

²⁴Ad oggi, la sola statistica di associazione di combinazioni di due elementi appare rigorosamente esplorata, sia per motivi di complessità matematica, sia per ragioni di semplicità pratica di implementazione degli algoritmi (molti processi, come quelli di clustering, hanno spesso alla base coppie di elementi) (Evert, 2004, p.32).

	w_2	$\neg w_2$			
w_1	O_{11}	O_{12}	$= R_1$	w_1	$E_{11} = \frac{R_1 C_1}{N}$ $E_{12} = \frac{R_1 C_2}{N}$
$\neg w_1$	O_{21}	O_{22}	$= R_2$	$\neg w_1$	$E_{21} = \frac{R_2 C_1}{N}$ $E_{22} = \frac{R_2 C_2}{N}$
	$= C_1$	$= C_2$	$= N$		

Figura 2.8: Forma generale di una tabella di contingenza per espressioni di due elementi (tratta da Evert 2008, p.1231). w_1 e w_2 rappresentano rispettivamente il primo ed il secondo componente; \neg indica la negazione. A sinistra è presente la tabella delle frequenze osservate: le quattro caselle includono le frequenze di cooccorrenza dei due elementi (O_{11}), le occorrenze del primo elemento con le parole diverse dal secondo (O_{12}), le occorrenze del secondo elemento con le parole diverse dal primo (O_{21}), le occorrenze in cui entrambi i componenti sono diversi dagli elementi considerati (O_{22}). Le quantità R e C a margine indicano le frequenze marginali per w_1 , w_2 , $\neg w_1$, $\neg w_2$; N è la grandezza del corpus in termini di token. A destra compare la tabella di contingenza per le frequenze aspettate (indicate con E), con le rispettive formule di calcolo.

cui *t-score* e *log-likelihood*), invece, pur considerando tale discrepanza, quantificano la magnitudine della probabilità contro l'ipotesi di indipendenza tra i componenti (significatività statistica).

Parallelamente alla definizione e allo studio delle differenti misure d'associazione, si è avuto un progressivo sviluppo di strumenti di estrazione che, grazie alle misure, fossero in grado di identificare automaticamente in un corpus le collocazioni empiriche presenti. Tra la fine degli anni '90 e i primi 2000, numerosi contributi²⁵ hanno aiutato a definire un insieme di metodologie applicabili con un buon grado di successo a corpora di diverso tipo e in diverse lingue (francese, tedesco, estone, oltre al sempre predominante inglese) per fini lessicografici e terminologici. In seguito, Evert (2004) sviluppa l'UCS toolkit, uno strumento costituito da un insieme di codici utili all'applicazione del vasto numero di misure di associazione a coppie di parole. Ma è solo con l'avvento di Sketch Engine²⁶ (Kilgariff *et al.*, 2004) e la sua interfaccia facilmente gestibile anche da chi non abbia conoscenze informatiche, che l'associazione statistica fra parole diviene un punto di partenza importante per ogni studio lessicografico, terminologico o di ricerca sui fenomeni multiparola²⁷. Lo strumento, infatti, è in grado di generare *word sketches*, vale a dire tabelle ricavate da un corpus²⁸ tramite l'utilizzo di diverse misure d'associazione (Lexical Compu-

²⁵Come quelli di Choueka (1988), Smadja (1993), Daille (1994, 1996), Justeson & Katz (1995), Kageura & Umino (1996), Lin (1998), Lemnitzer (1998), Nerima *et al.* (2003), Kaalep & Muischnek (2003).

²⁶<http://www.sketchengine.co.uk/>.

²⁷Sketch Engine viene usato per la prima volta nella composizione del Macmillan English Dictionary (Rundell, 2002), ma successivamente coinvolto in numerosi progetti lessicografici.

²⁸In Sketch Engine trova per la prima volta grande applicazione l'idea del *Web as a Corpus*

ting Ltd., 2014), che riassumono il comportamento grammaticale e collocazionale di una data parola. Sketch Engine sfrutta l'informazione sulla categoria e la relazione grammaticale precedentemente assegnata a tutti i *token* del corpus:

Rather than looking at an arbitrary window of text around the headword, we look, in turn, for each grammatical relation that the word participates in (Kilgariff *et al.*, 2004, p.107).

Il corpus necessita, cioè, oltre che del processo di *PoS-tagging*, anche del processo di *parsing* sintattico, e sarà chiaro, in questo caso, che il sistema si fonda su cooccorrenze sintattiche.

Sketch Engine deve il suo successo anche alla possibilità di operare analisi in un elevato numero di lingue (ad oggi 52²⁹) grazie alla vasta disponibilità di corpora immagazzinati. In questo modo esso fornisce anche una prova indiretta del fatto che, in generale, l'associazione non sia *language-dependent*, bensì che ogni lingua necessiti di pattern collocativi.

Ulteriori lavori successivi hanno continuato ad esplorare le potenzialità delle misure d'associazione nell'individuazione di polirematiche e collocazioni. In generale, però, si è assistito ad un progressivo spostamento delle metodologie dal raffinamento delle formule di calcolo statistico o dal loro uso *tout court*, verso l'integrazione della statistica con altri fattori linguistici, come la considerazione di pattern grammaticali, sintattici, ecc. (come già visto in Sketch Engine) al fine di una migliore individuazione automatica dei fenomeni multiparola. Il lavoro di Ramisch *et al.* (2010b), ad esempio, presenta un nuovo strumento in grado di assegnare automaticamente cinque misure d'associazione ad n-grammi con la possibilità di considerare le loro categorie grammaticali. I contributi di Villada Moirón & Tiedemann (2006), Ramisch *et al.* (2010a) e Pereira *et al.* (2014), esplorano la possibilità dell'uso combinato delle misure d'associazione e di corpora paralleli, sfruttando la proprietà, comune a molte espressioni idiomatiche, di non essere traducibili parola per parola da una lingua all'altra. Bannard (2007), Van de Cruys & Villada Moirón (2007), Ramisch *et al.* (2008) e Fazly *et al.* (2009) combinano l'informazione statistica con le tipiche caratteristiche di fissità sintattica e semantica delle *multiword expressions*. I lavori di Weller & Fritzinger (2010) e Cap *et al.* (2013) uniscono misure statistiche sia a informazioni derivanti da confronti interlinguistici, sia a comportamenti anomali morfosintattici, mentre infine Seretan (2011) presenta il primo lavoro incentrato sulla possibilità di utilizzare le misure d'associazione in combinazione con pattern sintattici derivanti dal *parsing* per l'estrazione automatica di fenomeni multiparola.

(Kilgariff & Grefenstette, 2001), tramite l'utilizzo di corpora di grandi dimensioni costituiti da numerosi testi collezionati in rete.

²⁹ottobre 2014, nda.

2.4.4 Valutazione dei risultati

In generale, tanto le singole misure d'associazione quanto gli approcci "ibridi" visti nel precedente paragrafo si dimostrano utili all'estrazione automatica di fenomeni multiparola dai testi. Le loro performance dipendono, però, da diversi fattori, quali le caratteristiche del corpus, la tipologia delle espressioni considerate, la lingua in esame.

Una volta stabilita, ad esempio, una soglia di significatività (relativamente al punteggio d'associazione) al di sopra della quale tutte le espressioni vengono incluse nell'insieme delle collocazioni empiriche, si è soliti riscontrare in tale insieme anche *errori*, cioè espressioni selezionate che non sono, in realtà, polirematiche o collocazioni. Queste ultime vengono definite, in gergo, *false positives* (fp), in opposizione alle espressioni correttamente individuate, dette *true positives* (tp).

Un punto critico, tuttavia, rimane la determinazione della natura delle espressioni. Come ricorda Evert (2004, p. 26): «Most approaches assume that there is a *binary distinction* between collocational and non-collocational pairs», cosa che risulta completamente infondata in virtù del continuum su cui polirematiche, collocazioni ed espressioni libere si muovono. Pur con la consapevolezza del limite di tale metodologia, una delle possibilità nell'analisi dei *true and false positives* è quella di confrontare le espressioni identificate rispetto ad una lista di espressioni esterna considerata affidabile e composta da un elevato numero di espressioni dallo status categoriale definito: un cosiddetto *gold standard*. Specie per le analisi di terminologia settoriale, è possibile ipotizzare che il gold standard rappresenti una buona approssimazione di tutte le espressioni multiparola del linguaggio tecnico-specialistico considerato. In questo caso, il riconoscimento di un'espressione identificata dalla misura d'associazione come una reale espressione multiparola è subordinata alla presenza di quest'ultima nel *gold standard*. In altri casi (i.a. Lapata *et al.* 1999), una volta ottenuta la lista di espressioni selezionate dalla metodologia in esame, è possibile chiedere a degli annotatori umani (ad esempio linguisti o parlanti della lingua in oggetto) di etichettare le espressioni come fenomeni multiparola o espressioni libere. Successive statistiche sull'eventuale accordo o disaccordo tra gli annotatori potranno produrre un giudizio definitivo su ogni espressione.

In ultima istanza, l'identificazione di ciò che deve o non deve essere incluso nell'insieme delle espressioni multiparola dipende dagli interessi e obiettivi della ricerca. Krenn *et al.* (2004) forniscono la seguente sintesi:

The phenomena subsumed by lexical collocations are manifold, ranging from lexical proximities in texts to syntactic and semantic units showing semantic opacity, and syntactic irregularity and rigidity. Accordingly there is a variety of definitions and terminology. Both are influenced by different linguistic traditions and by the particular computational linguistics applications for which collocations are considered to be relevant. *We typically find an opportunistic approach to collocativity, i.e., the definition of TPs [true positives] depends on the intended application*

rather than being motivated by (linguistic) theory, and it covers a mixture of different phenomena and classes of collocations. Moreover, even when well-defined criteria and explicit annotation guidelines are available, annotators may make different decisions for some of the collocation candidates, because of mistakes, differences in their intuition, different interpretation of the guidelines, etc. All this makes it hard to give a systematic experimental account of the true usefulness of a certain AM [association measure] for collocation extraction (Krenn *et al.*, 2004, p. 1, c.vo mio).

Gli stessi, tuttavia, dimostrano che annotazioni fornite da sole persone con elevata «linguistic expert knowledge» (*ibid.*, p. 8) risultano l'unico modo per ottenere, a livello pratico, un buon accordo e una valutazione affidabile.

2.4.5 Limiti delle misure d'associazione

Si è detto, nel paragrafo 2.3, come esistano delle preferenze di combinazione tra *item* lessicali, il cui status può essere controverso. Un *libro*, ad esempio, può essere più comunemente *letto*, *sfogliato*, *aperto*, *chiuso*, *riposto*, *prestato*, anche senza che *leggere* o *chiudere un libro* vengano generalmente definite collocazioni. In questi casi, l'eventuale comparsa di tali espressioni nell'insieme delle collocazioni empiriche in analisi attesterà la possibilità delle misure d'associazione o di procedure che le coinvolgono nel processo di identificazione di fenomeni multiparola di selezionare *false positives*.

Allo stesso tempo si è visto come la collocatività empirica possa essere considerata un epifenomeno motivato da cause di carattere fraseologico: «idioms, lexical collocations, clichés, cultural stereotypes, semantic compatibility and many other factors are hidden causes that result in the observed associations between words» (Evert, 2008, p.1218). Tali associazioni, tuttavia, quantificate dal punteggio d'associazione generato da un determinato processo, non sembrano disporre le collocazioni empiriche su un asse che riproduca la polarizzazione intuitiva tra espressioni rappresentanti unità semantiche e cristallizzate (polirematiche) ed espressioni preferenziali (collocazioni lessicali). Non sembra esistere, cioè, una forte correlazione tra il maggiore punteggio di associazione statistica e il maggior grado di idiomaticità, fissità, unità di un'espressione. In altre parole, le misure d'associazione non riescono a discriminare le espressioni multiparola sulla base della loro diversa natura, poiché non ricavabile da mere informazioni di cooccorrenza. A titolo di esempio si riporta la tabella di Figura 2.9, in cui sono mostrate coppie di parole estratte dal British National Corpus (BNC, 2001) ed ordinate secondo il punteggio di *log-likelihood*.

Come si vede, espressioni del tipo *ask (the) secretary* o *write (a) letter*, che potremmo considerare libere, appaiono mescolate e con punteggi maggiori di espressioni più tipicamente collocative, come *solve (a) problem* o *raise (a) question* o polirematiche, come *make sense*.

word pair	freq.	association	word pair	freq.	association
<i>take place</i>	7606	41942.15	<i>meet needs</i>	520	4183.22
<i>play role</i>	1488	11710.46	<i>make mistake</i>	763	4114.63
<i>open door</i>	1438	11299.73	<i>make decision</i>	1172	3943.51
<i>see chapter</i>	1461	9795.36	<i>keep eye</i>	577	3671.20
<i>give rise</i>	1499	9521.99	<i>tell storey</i>	527	3616.61
<i>make sense</i>	1888	7996.27	<i>show sign</i>	533	3577.90
<i>take advantage</i>	1557	7529.19	<i>pay tribute</i>	336	3390.79
<i>see page</i>	1294	7374.60	<i>thank goodness</i>	224	3338.03
<i>play part</i>	1331	7359.75	<i>take action</i>	1023	3302.98
<i>draw attention</i>	836	6610.98	<i>shake hand</i>	342	3289.48
<i>answer question</i>	743	6558.03	<i>take step</i>	759	3271.63
<i>take part</i>	2358	6424.23	<i>get hold</i>	614	3265.61
<i>ask question</i>	898	6373.39	<i>form basis</i>	448	3191.22
<i>take care</i>	1295	6196.21	<i>ring bell</i>	235	3093.21
<i>ask secretary</i>	621	5932.39	<i>closed door</i>	346	3091.96
<i>solve problem</i>	645	5706.53	<i>shut door</i>	322	3039.70
<i>wait minute</i>	428	5422.05	<i>write letter</i>	445	3023.47
<i>make use</i>	1441	4954.34	<i>give impression</i>	638	2948.46
<i>take account</i>	1164	4626.66	<i>make contribution</i>	682	2890.16
<i>form part</i>	886	4335.17	<i>raise question</i>	555	2882.07

Figura 2.9: Coppie di elementi del pattern verbo + sostantivo (oggetto) estratti dal British National Corpus (BNC) in ordine decrescente secondo il punteggio d'associazione di log-likelihood (tratta da Evert 2004, p.22).

Considerazioni di questo tipo mostrano come la statistica possa porsi come obiettivo l'individuazione di espressioni in cui i componenti condividano un legame empiricamente 'particolare' ma non riesca a dirimere il continuum categoriale dei fenomeni collocativi.

Si mostrerà, invece, come lo studio quantitativo del comportamento sintattico e semantico delle espressioni frequenti possa contribuire a porre le basi di una differenziazione empirica del vasto insieme dei fenomeni multiparola.

3

Uno strumento per le analisi variazionali dei fenomeni multiparola

3.1 Approcci computazionali e categorizzazione

Come accennato nel precedente capitolo, in generale i moderni approcci ai fenomeni multiparola che coinvolgono metodologie computazionali hanno privilegiato strategie d'indagine volte alla loro estrazione piuttosto che a riflessioni teoriche che potessero, in particolare, dividere il loro grande *mare magnum* di variabilità in categorie. In altri termini, gran parte della ricerca linguistico-statistica si è orientata verso l'individuazione di tecniche che separassero l'insieme delle *multiword*, considerate come un tutt'uno, dalle espressioni libere, relegando i problemi correlati alle sfumature di confine tra i due gruppi alle definizioni stabilite a monte delle analisi, generalmente variabili a seconda degli studi.

Pochi sono stati, invece, i contributi che, pur considerando la generale non nitidezza dei confini categoriali nei fatti di lingua, abbiano indagato in senso computazionale e statistico eventuali raggruppamenti interni allo stesso insieme di collocazioni empiriche individuato. In quest'ottica, riuscire a dirimere con procedure automatiche o semiautomatiche l'insieme dei fenomeni multiparola avrebbe un duplice effetto di grande importanza. Da un lato, sul versante della teoria linguistica, implicherebbe il successo nell'individuazione di principi, tendenze o regolarità generali alla base della diversa natura delle differenti tipologie di espressioni multiparola (prime fra tutte polirematiche e collocazioni). Dall'altro, una metodologia automatica di differenziazione avrebbe grandi ricadute applicative, specie in lessicografia: sarebbe possibile, ad esempio, a valle dell'estrazione da un corpus, avere distinte le espressioni idiomatiche che necessitano di inclusioni specifiche nel lessico della lingua, dalle collocazioni che potrebbero essere incluse in un dizionario combinatorio, da combinazioni istituzionalizzate, che potrebbero essere utilizzate come esempi nelle definizioni dei lemmi di un dizionario generale, e così via.

In virtù di quanto esposto nel par. 2.4.5, risulta chiaro che la base per lo sviluppo di procedure in grado di indagare più nello specifico eventuali categorizzazioni inter-

ne all'insieme dei fenomeni multiparola trascende la mera statistica. Storicamente, nonostante l'obiettivo di molti studi fosse ancora una individuazione automatica *tout court* di ciò che fosse o meno *multiword*, alcuni lavori che hanno provato ad integrare l'informazione d'associazione statistica a misure di motivazione linguistica hanno posto le basi per i futuri sviluppi sulla categorizzazione.

Il lavoro di Lin (1999), ad esempio, propone di utilizzare un test linguistico di sostituibilità per l'individuazione automatica di *multiword expressions* inglesi, sfruttando il concetto di non composizionalità su una scala graduale. L'idea di base è che più un'espressione è non composizionale, più la sua misura d'associazione statistica (*mutual information* nello studio in esame) differirà da quella calcolata per la stessa espressione in cui uno dei componenti è stato sostituito da una parola ad esso simile¹. I risultati sono promettenti, benché, per ammissione dell'autore, viziati da errori sistematici.

Nei primi anni 2000, alcuni contributi come quelli di Baldwin *et al.* (2003), Bannard *et al.* (2003) e McCarthy *et al.* (2003) hanno continuato ad esplorare la possibilità di identificare attraverso metodologie computazionali una scala continua e ordinata di espressioni multiparola sulla base della sola variabile della non composizionalità, focalizzandosi sempre sulla lingua inglese. Ciò che accomuna i suddetti lavori è il tentativo di determinare diversi livelli di composizionalità in base a misure di *similarità distribuzionale* (cfr. par. 3.3.3.2), ovvero basandosi sull'idea che il significato di due generiche entità lessicali è tanto più simile quanto più esse occorrono (si *distribuiscono*) in cotesti comuni.

Baldwin *et al.* (2003) utilizzano tecniche di *latent semantic analysis*² per confrontare la similarità semantica tra l'intera espressione multiparola immersa in contesto e i suoi componenti, studiando *phrasal verbs* (*take off*) e composti NN (*motor car*), affermando che maggiore è la similarità tra espressione e i componenti presi singolarmente, maggiore risulterà la composizionalità dell'espressione.

Bannard *et al.* (2003), invece, si soffermano solo sulle espressioni verbali, provando a quantificare con diverse tecniche statistiche il contributo della preposizione del *phrasal verb* nella sua accezione indipendente, al significato dell'espressione in cui essa appare, ottenendo risultati soddisfacenti nell'ordinare le espressioni tra i due opposti poli di completa opacità semantica o completa trasparenza.

McCarthy *et al.* (2003), infine, utilizzano sempre sui *phrasal verbs* una serie di metodi che coinvolgono la sostituzione dell'espressione o di parti di essa con un certo numero di *nearest neighbours*, vale a dire sinonimi identificati empiricamente come occorrenti negli stessi cotesti. Attraverso la quantificazione della distribuzionalità dell'espressione originale e di quella sostituita da sinonimi, si assegna un punteggio

¹Nello studio la similarità è intesa sia in un senso sinonimico (*economic effect/consequence/impact*) che di appartenenza ad una stessa classe di potenziale sostituibilità (*red/yellow/black tape*).

²La *latent semantic analysis* (LSA) è una tecnica che astrae concetti o polarizzazioni che riguardano gli elementi di un testo e i frammenti di testo in cui essi occorrono sulla base di relazioni di cooccorrenza. Per approfondire si rimanda a Dumais (2005) o a Bolasco (1999) per l'italiano.

di composizionalità, che risulta correlato al grado di composizionalità assegnato alle espressioni da annotatori umani. Questo lavoro sembra essere il primo a notare un aspetto potenzialmente fondamentale nella gestione computazionale delle espressioni multiparola con metodologie diverse dalle classiche misure d'associazione. A fronte di risultati distribuzionali in miglior accordo con i giudizi umani rispetto all'uso delle tecniche puramente statistiche, gli autori affermano che «it might be better not to filter with statistics before looking at compositionality using an automatic thesaurus» (McCarthy *et al.*, 2003, p. 79). Esiste, cioè, la possibilità che criteri squisitamente linguistici, una volta automatizzati e applicati a grandi moli di dati per ottenerne riscontri empirici, siano in grado di cogliere ciò che sfugge alle misure d'associazione.

Un'ulteriore e importante riflessione, presente sia nel lavoro di McCarthy e colleghi che in quello del gruppo di Bannard, è che la composizionalità, come c'era da aspettarsi, non riesce ad essere l'unico asse di oscillazione delle diverse espressioni assimilabili ai fenomeni multiparola e per questo ne limita lo spettro categoriale. Il fatto che sia permessa la sostituzione di un componente con un lemma ad esso semanticamente vicino può essere un buon indicatore di produttività dell'espressione, e per questo della sua natura composizionale; tuttavia, l'impossibilità di sostituire i componenti non risulta, in maniera complementare, un forte indice di non composizionalità. Come precisano McCarthy *et al.* (2003, p. 75): «an institutionalised non-productive combination, such as *frying pan*, would not have near neighbour substitutes, but would nevertheless be compositional».

Nonostante queste riflessioni, un successivo lavoro di Venkatapathy & Joshi (2005) si sofferma ancora sullo studio della composizionalità delle espressioni multiparola, focalizzandosi sulle espressioni VN, proponendo una metodologia che rappresenta una *summa* delle strategie sviluppate all'epoca. Gli autori espongono, infatti, un approccio che generi un ordinamento delle espressioni grazie a un punteggio ottenuto dall'integrazione di diversi indici. Sul versante statistico vengono presi in considerazione la frequenza dell'espressione e una misura d'associazione (*mutual information*). Vengono quindi integrate le informazioni ottenute grazie a test di sostituibilità mutuati dall'approccio di Lin (1999), test sulla frequenza di distribuzione dell'oggetto (Tapanainen *et al.*, 1998), secondo cui se un certo oggetto appare solo con pochi verbi, le espressioni in cui è presente tendono ad essere di natura idiomatica, e test distribuzionali basati sulla *latent semantic analysis*. Anche qui i risultati mostrano come l'uso integrato dei diversi criteri linguistici produca un accordo maggiore delle misure d'associazione rispetto alle annotazioni dei valutatori umani per quanto riguarda il *ranking* delle espressioni su una scala di composizionalità³.

³Nelle linee guida per l'annotazione gli autori definiscono una serie di livelli in cui collocare le espressioni. Come spesso accade nelle definizioni specificamente fornite a seconda degli studi, i criteri guida risultano discutibili. L'espressione *leave a mark* viene inclusa nel livello di massima non composizionalità, poiché «no word in the expression has any relation to the actual meaning of the expression» (Venkatapathy & Joshi, 2005, p. 901), nonostante *mark* viene adoperato in uno dei suoi principali significati metaforici. Un altro livello è definito dalle espressioni sostituibili da

Le riflessioni di Bannard, McCarthy e colleghi sulla limitatezza dello studio della sola non composizionalità sembrano essere colte in un lavoro quasi contemporaneo di Wermter & Hahn (2004), in cui viene proposto un criterio di individuazione dei fenomeni multiparola diverso dalla composizionalità e basato su una misura che quantifichi statisticamente la preferenzialità di modificazione dell'espressione per inserzione. Gli autori, infatti, studiano la possibilità per le combinazioni tedesche PNV (es. *zur Verfügung stehen*) di essere interrotte da altro materiale linguistico. L'assunto di base è l'ipotesi che combinazioni di questo tipo siano meno modificabili, e per questo identificabili come collocazioni, se ammettono un particolare elemento di inserzione tra i componenti predominante rispetto a qualsiasi altro attestato (Wermter & Hahn, 2004, p. 982), come nel caso di *unter Druck geraten* dove il modificatore preferenziale di *Druck* risulta *politischen*. Grazie a un confronto sull'identificazione automatica di collocazioni tra la nuova misura proposta e le già note misure puramente statistiche, gli autori confermano che prendere in considerazione un criterio linguistico genera performance migliori. Gli autori, inoltre, compiono un'importante verifica, mostrando empiricamente che la preferenzialità di un qualche elemento specifico inserito tra i componenti rispetto alla scelta libera è un tratto tipico delle collocazioni rispetto alle espressioni definite standard, che invece mostrano una più regolare equidistribuzione delle parole intervenienti. Un tale risultato, tacitamente e intuitivamente dato per assodato in molti studi teorici, viene così per la prima volta supportato da prove oggettive.

Nonostante, tuttavia, nel lavoro di Wermter & Hahn venga introdotto il nuovo elemento della modificabilità, l'obiettivo dello studio rimane quello della mera individuazione e separazione delle espressioni multiparola dall'insieme delle espressioni libere.

Fazly & Stevenson (2007) sono i primi a sviluppare una procedura il cui obiettivo sia esplicitamente quello dell'identificazione automatica di diverse classi di fenomeni multiparola, focalizzandosi sulle combinazioni inglesi di verbo e oggetto. Gli autori riconoscono *a priori* l'esistenza di quattro classi distinte relative al pattern VN preso in considerazione: *idioms* (fortemente idiosincratici e non composizionali, come ad esempio *shoot the breeze*), costruzioni a verbo supporto (in cui il verbo è svuotato del proprio significato e investito di quello dell'oggetto, come nel caso di *make a decision*), costruzioni astratte (combinazioni frequenti o istituzionalizzate in cui il verbo è generalmente utilizzato metaforicamente in un suo senso astratto, come ad esempio *make a living*), combinazioni libere (*drink a coffee*). Analogamente a Venkatapathy & Joshi (2005), gli autori costruiscono una serie di indici, ognuno volto a individuare una determinata caratteristica idiosincratica. L'istituzionalizzazione dell'espressione viene individuata anche qui dalla misura d'associazione *mutual information* ed essa è supposta essere una caratteristica saliente di tutte le espressioni eccetto quelle libere. Vengono poi prese in considerazione le fissità lessicale e sintat-

un singolo verbo (come *take a look*), ma si è già visto nel par. 1.5.2 quanto il criterio della parafrasi risulti problematico.

tica (pregnanti per le sole classi di *idioms* e verbi-supporto): la prima è valutata da un indice che calcola la vicinanza statistica (*z-score*) tra il punteggio d'associazione dell'espressione e la media di quelli ottenuti da sue varianti in cui uno dei componenti è stato sostituito; la seconda tramite una misura dell'entropia delle diverse varianti sintattiche dell'espressione rispetto al comportamento tipico di una combinazione verbo oggetto. Infine la non composizionalità è valutata in base al confronto distribuzionale del contesto dell'espressione rispetto a quelli dei suoi singoli componenti, ed essa è considerata una caratteristica tipica dei soli *idioms* e in misura minore delle combinazioni a verbo-supporto.

Un sistema di alberi decisionali è in grado di raggruppare, sulla base dei valori di ciascun indice, le espressioni estratte da un corpus nelle quattro categorie ipotizzate. Fazly & Stevenson mostrano che le performance di classificazione ottenute dalla procedura automatica sono variabili a seconda dell'indice considerato. Ad esempio, nel computo dell'indice di composizionalità, l'eliminazione del confronto tra la similarità di contesto dell'espressione e quello del solo verbo migliora i risultati, in quanto i verbi, specie nelle costruzioni a verbo supporto, presentano grande polisemia e quindi un significato non unicamente definito. Nonostante l'istituzionalizzazione misurata attraverso la misura d'associazione ottenga discreti risultati, gli indici di fissità lessicosintattica raggiungono l'accuratezza maggiore nell'accordo di classificazione con gli annotatori.

I lavori qui citati si configurano, quindi, come studi che accrescono in misura sempre maggiore la rilevanza della ricerca empirica sul comportamento delle espressioni multiparola al fine di individuarle e categorizzarle. Va notato che a seconda delle ipotesi (intuitive) alla base della ricerca, gli studi si concentrano su specifiche caratteristiche, ma soprattutto su un gruppo limitato di espressioni che, ad eccezione di una parte del lavoro di Baldwin *et al.* (2003), sono costituite da sole espressioni verbali. Manca, quindi, un'indagine più organica che abbracci altre categorie di espressioni e che metta in relazione le eventuali differenze di classificazione rispetto ai verbi.

3.2 Motivazioni per uno studio della modificabilità empirica

Il quadro che emerge dagli studi del primo decennio del secolo citati in precedenza può essere riassunto dalle seguenti indicazioni, seppur parziali e relative a gruppi ristretti di espressioni:

- prendere in considerazione elementi di natura puramente linguistica nella metodologia di analisi migliora le performance di estrazione delle espressioni multiparola;
- le espressioni multiparola hanno un comportamento idiosincratico empiricamente verificabile rispetto alle espressioni che consideriamo 'libere' (Werm-

ter & Hahn, 2004), le quali hanno tendenzialmente una maggiore libertà di modificazione;

- la fissità lessicosintattica, vale a dire l'impossibilità di modificare sintagmaticamente o paradigmaticamente i componenti dell'espressione è una delle caratteristiche più rilevanti ai fini di una categorizzazione (Fazly & Stevenson, 2007).

Risulta naturale chiedersi, a questo punto, se uno studio a più ampia copertura del comportamento variazionale di polirematiche e collocazioni in merito alla modificabilità morfo-sintattico-semantiche possa aiutare a definire una migliore categorizzazione generale che valga per diverse tipologie di espressioni multiparola e non solo per i verbi. Se, infatti, le misure d'associazione hanno stabilmente dimostrato la propria efficacia (a diversi gradi) nell'individuazione di tali espressioni, è necessario mettere alla prova le potenzialità dei test strettamente linguistici anche e soprattutto quando non si abbia un loro abbinamento con metodologie statistiche, al fine di evidenziare il loro singolo contributo. Se l'obiettivo resta la pura categorizzazione, infatti, non è necessario vedere la statistica come un passaggio preliminare necessario all'individuazione di un nucleo iniziale di candidati sostanzialmente *multiword*, bensì sono le stesse variabilità a poter categorizzare il continuum lessico-sintattico dalle espressioni fisse a quelle standard.

In quest'ottica si può procedere, quindi, a un ribaltamento della prospettiva, tradizionalmente istituzionalizzata fino al presente, di utilizzare il contributo dato dalle caratteristiche idiosincratiche di modificabilità lessicosintattica al fine di migliorare l'individuazione automatica statistica di *multiword* (come nel recente lavoro di Cap *et al.* 2013), per passare a una visione in cui l'obiettivo primario è la categorizzazione attraverso le suddette modificabilità (seguendo la via indicata da Fazly & Stevenson 2007) e, solo di conseguenza, la separazione delle espressioni multiparola da quelle libere.

È però interessante anche l'idea di ottenere una categorizzazione delle espressioni multiparola senza una definizione *a priori* delle sottoclassi in cui collocarle (come invece fatto in Fazly & Stevenson 2007), lasciando che siano i dati ottenuti dallo studio delle modificabilità a guidare la definizione dei gruppi in cui le espressioni risultano distribuite. Nel presente lavoro si è quindi optato per quest'ultima scelta.

3.3 Metodologia di studio delle variabilità

Al fine di operare uno studio sulle possibili modificazioni delle espressioni multiparola, si è scelto di costruire uno strumento computazionale *ad hoc* che fosse in grado di effettuare una serie di test in maniera automatica su un vasto insieme di espressioni, come esplicitato nel seguito. Il carattere **empirico** dello strumento è giustificato da due principali ragioni: in primis, le espressioni di cui verrà studiata la modificabilità sono estratte da un corpus che costituisce la base empirica dello studio; in secondo luogo, la variabilità lessicosintattica delle espressioni verrà valutata

in base all'eventuale attestazione delle loro varianti nello stesso corpus. Il risultato è, inoltre, **quantitativo** poiché per ognuna delle espressioni prese in esame, grazie al programma, sarà possibile quantificare le occorrenze nel corpus di ognuna delle varianti attestate e tradurre tale dato in un indice percentuale di modificabilità.

3.3.1 Corpus, pattern e lista di espressioni

Il corpus, che si configura come l'elemento essenziale dell'intero processo al fine dell'utilizzo del *tool*, necessita dell'etichettatura *part-of-speech* per ciascuno dei suoi *token* che devono anche essere stati ricondotti al proprio lemma di riferimento. Il *PoS-tagging* è un requisito fondamentale in quanto i test effettuati variano a seconda del pattern grammaticale specificato in input, come precisato nel seguito.

Qualora il corpus abbia subito anche un processo di *parsing* sintattico, tale livello di informazione risulterà utile (ma non strettamente necessario) all'analisi, poiché esso permetterà di includere nei risultati dei test espressioni i cui componenti abbiano maggiore flessibilità di movimento, ferma restando la corretta individuazione dei legami tra di essi⁴.

Dal corpus è necessario poi estrarre la lista di espressioni di cui andrà testata la modificabilità. Una scelta possibile sarebbe quella di considerare una lista di candidati ad espressioni multiparola selezionata grazie ad una misura d'associazione, tuttavia tale decisione potrebbe introdurre un importante *bias* nei risultati, a causa della possibile tendenza della misura ad individuare solo particolari tipi di espressioni multiparola (Evert & Krenn, 2001). Va sottolineato, infatti, che lo scopo dello studio è individuare una categorizzazione delle espressioni in base al loro differente comportamento variazionale e va quindi garantita la neutralità della base empirica che determinerà la categorizzazione. In ragione di ciò si è preferito considerare, per ogni pattern grammaticale analizzato, l'insieme delle espressioni più frequenti che soddisfino lo stesso pattern, in un numero che può essere specificato in input, a seconda dell'esigenza di ricerca o delle dimensioni del corpus. In questo modo verrà selezionato un insieme di espressioni contenente, verosimilmente, tanto espressioni multiparola che combinazioni libere i cui componenti hanno alta frequenza (cfr. Par. 2.4.2). La scelta di considerare le espressioni con un alto numero di occorrenze è motivata anche dal fatto che i test esposti di seguito si basano sul concetto che, data un'espressione, essa è tanto più modificabile quanto più un certo tipo di modifica è attestata nel corpus (e, analogamente, la non attestazione nel corpus di varianti dell'espressione implica l'impossibilità di quest'ultima ad essere modificata), legando a filo stretto la realtà empirica al piano teorico. L'affidabilità delle conclusioni traibili da premesse di questo tipo dipende, oltre che dal corpus stesso (supposto di dimensioni adeguate e rappresentativo), dal numero di occorrenze delle espressioni da testare: un alto numero di occorrenze garantisce, infatti, un'ampia base su cui

⁴Un test che verifichi, ad esempio, l'interrompibilità dell'espressione potrà quantificare le possibili varianti senza fissare uno *span* determinato di parole intervenienti.

quantificare in maniera verosimile le proporzioni della modificabilità del fenomeno, avvicinando il risultato empirico a quello teorico.

L'insieme di partenza delle espressioni da analizzare attraverso lo strumento computazionale comprende quindi le n espressioni più frequenti del corpus, rispondenti a un determinato pattern, e considerate nella loro *forma base lemmatizzata*, in cui cioè il numero di occorrenze è stato calcolato considerando i lemmi di ciascun componente nella combinazione identificata dal pattern.

Motivazioni legate al carico computazionale di calcolo ed esecuzione delle analisi hanno determinato, per semplicità, lo sviluppo dei test per la gestione di pattern grammaticali di due o tre elementi, contenenti, in quest'ultimo caso, al massimo due parole piene.

Va precisato, infine, che tutte le conclusioni traibili dalle analisi su un tale campione hanno l'ambizione di dar conto del comportamento generale delle espressioni relative ad un determinato pattern in virtù della loro affidabilità in quanto altamente frequenti. Tuttavia, non è da escludere che espressioni di bassa frequenza potrebbero esibire comportamenti e quindi categorizzazioni diversi. Tale limite metodologico non è però superabile, in quanto non risulterebbe sensato studiare le proporzioni di modificabilità attestata e quindi giudicare se un'espressione sia modificabile o meno sulla base di un numero di frasi in cui essa appare nel corpus dell'ordine delle unità.

3.3.2 Variazioni sintagmatiche

Per ognuna delle espressioni che compongono la lista iniziale in input, lo strumento è in grado di compiere una serie di test sulla variabilità sintagmatica, vale a dire la possibilità di muovere i costituenti o interromperli con altro materiale linguistico. In generale la valutazione della modificabilità è ottenuta sulla base della quantificazione del numero di occorrenze delle varianti rispetto al numero di occorrenze dell'espressione nella sua forma base, quest'ultimo definito d'ora in avanti n_{bf} . In tutti i casi le ricerche operate nel corpus (*query*) vengono effettuate considerando i lemmi dei componenti.

Per chiarezza di notazione, nell'intero paragrafo w_i rappresenterà il generico componente dell'espressione, con $i = 1, 2$ nel caso di pattern a due elementi o $i = 1, 2, 3$ nel caso di pattern a tre elementi. La sequenza di componenti con pedici in ordine crescente rappresenta la forma base. Con il simbolo [] si indicherà, invece, materiale linguistico interveniente, costituito da almeno un token.

3.3.2.1 Interrompibilità

Il primo dei test che lo strumento è in grado di effettuare è il test di interrompibilità, in cui si ricercano nel corpus eventuali attestazioni dell'espressione considerata in cui i componenti siano separati da altre parole. Nel caso di pattern a due elementi (es. NA \rightarrow *punto debole*), indicata con $w_1 w_2$ la forma base dell'espressione,

lo schema generico di query per la ricerca di varianti interrotte è il seguente:

$$\begin{array}{ccc} w_1 & [] & w_2 \\ N & [] & A \\ \text{punto} & \text{più} & \text{debole} \end{array} \quad (3.1)$$

Nel caso di pattern a tre elementi (es. $\text{VDN} \rightarrow \text{fare la spesa}$), detta w_1 w_2 w_3 la forma base, le query effettuate risultano:

$$\begin{array}{ccc} w_1 & [] & w_2 & w_3 \\ V & [] & D & N \\ \text{farò} & \text{domani} & \text{la} & \text{spesa} \end{array} \quad (3.2)$$

$$\begin{array}{ccc} w_1 & w_2 & [] & w_3 \\ V & D & [] & N \\ \text{fatta} & \text{la} & \text{grande} & \text{spesa} \end{array}$$

Nel caso in cui l'unico livello di annotazione disponibile sia il PoS-tagging, è possibile stabilire lo *span* di parole intervenienti. Se, invece, il corpus ha subito un processo di parsing sintattico, lo *span* può essere arbitrario, purché limitato dal confine della frase, poiché i costituenti risultano collegati indipendentemente dal numero di parole intervenienti.

Detto n_{int} il numero totale di occorrenze delle varianti dell'espressione attestate nel corpus che presentano interruzione, l'indice di interrompibilità I_{syn}^{int} è dato dalla seguente formula:

$$I_{syn}^{int} = \frac{n_{int}}{n_{bf} + n_{int}} \quad (3.3)$$

In questo modo l'indice misura la proporzione delle espressioni interrotte rispetto al totale delle espressioni (interrotte e nella forma base), aumentando all'aumentare delle prime, ma rimanendo vincolato ad assumere valori tra 0 e 1.

3.3.2.2 Ordine modificabile

Un ulteriore test riguarda la possibilità di modificare l'ordine di occorrenza delle parole piene componenti l'espressione. Nel caso di pattern a due elementi (es. $\text{NA} \rightarrow \text{infrarosso lontano}$), lo schema di query è il seguente:

$$\begin{array}{cc} w_2 & w_1 \\ A & N \\ \text{lontano} & \text{infrarosso} \end{array} \quad (3.4)$$

Nel caso di pattern a tre elementi (es. $\text{NCN} \rightarrow \text{flora e fauna}$), l'inversione riguarda il primo e l'ultimo componente:

$$\begin{array}{ccc} w_3 & w_2 & w_1 \\ N & C & N \\ \text{fauna} & \text{e} & \text{flora} \end{array} \quad (3.5)$$

Analogamente all'indice di interrompibilità, l'indice relativo all'ordine modificabile I_{syn}^{ord} è calcolato dalla seguente formula, dove n_{ord} è il numero di occorrenze dell'espressione con i componenti invertiti:

$$I_{syn}^{ord} = \frac{n_{ord}}{n_{bf} + n_{ord}} \quad (3.6)$$

Anche in questo caso l'indice valuta la proporzione di modifiche attestate rispetto all'insieme delle espressioni modificate e in forma base, muovendosi nell'intervallo delimitato da 0 e 1.

3.3.2.3 Trasformazioni sintattiche

Per i soli pattern verbali (es. $VDN \rightarrow \textit{disputare il campionato}$) lo strumento è in grado di effettuare una serie di test per individuare nel corpus le attestazioni dell'espressione in una delle sue possibili trasformazioni sintattiche.

Il test di **topicalizzazione** ricerca nel corpus l'espressione in cui il complemento verbale (diretto o indiretto) è anteposto senza riprese successive al verbo, quest'ultimo flessa anche nelle sue forme analitiche:

$$\begin{array}{ccccc} w_2 & w_3 & & w_1 & \\ D & N & ([\]*) & V & \\ \text{il} & \text{campionato} & (\text{io}) & \text{disputo} & \end{array} \quad (3.7)$$

$$\begin{array}{ccccc} w_2 & w_3 & & w_1 & \\ D & N & ([\]*) & V_{aux} & V_{pp} \\ \text{il} & \text{campionato} & (\text{io}) & \text{ho} & \text{disputato} \end{array} \quad (3.8)$$

Il test di **ripresa anaforica** ricerca eventuali attestazioni dell'espressione verbale in una delle possibili costruzioni con ripresa anaforica pronominale dell'oggetto, secondo i seguenti schemi:

$$\begin{array}{ccccc} w_2 & w_3 & & w_1 & \\ D & N & ([\]*) & Pron & V \\ \text{il} & \text{campionato} & (\text{io}) & \text{lo} & \text{disputo} \end{array} \quad (3.9)$$

$$\begin{array}{ccccc} w_2 & w_3 & & w_1 & \\ D & N & ([\]*) & Pron & V_{aux} & V_{pp} \\ \text{il} & \text{campionato} & (\text{io}) & \text{l'} & \text{ho} & \text{disputato} \end{array} \quad (3.10)$$

Il test di **passivizzazione** ricerca eventuali attestazioni di costruzioni passive dell'espressione in esame secondo i seguenti schemi:

$$\begin{array}{ccccc} w_2 & w_3 & & w_1 & \\ D & N & & V_{aux} & V_{pp} \\ \text{il} & \text{campionato} & & \text{viene} & \text{disputato} \end{array} \quad (3.11)$$

$$\begin{array}{ccccc}
w_2 & w_3 & & & w_1 \\
D & N & V_{aux} & V_{aux} & V_{pp} \\
\text{il} & \text{campionato} & \text{è} & \text{stato} & \text{disputato}
\end{array} \quad (3.12)$$

Il test di **relativizzazione**, infine, quantifica la presenza di occorrenze dell'espressione in cui l'oggetto viene ripreso da un pronome relativo ad introduzione di una proposizione subordinata, secondo i seguenti schemi:

$$\begin{array}{ccccc}
w_2 & w_3 & & & w_1 \\
D & N & Pron & ([\]*) & V_{pp} \\
\text{il} & \text{campionato} & \text{che} & (\text{io}) & \text{disputo}
\end{array} \quad (3.13)$$

$$\begin{array}{ccccc}
w_2 & w_3 & & & w_1 \\
D & N & Pron & ([\]*) & V_{aux} & V_{pp} \\
\text{il} & \text{campionato} & \text{che} & (\text{Maria}) & \text{ha} & \text{disputato}
\end{array} \quad (3.14)$$

$$\begin{array}{ccccc}
w_2 & w_3 & & & w_1 \\
D & N & Pron & V_{aux} & V_{aux} & V_{pp} \\
\text{il} & \text{campionato} & \text{che} & \text{è} & \text{stato} & \text{disputato}
\end{array} \quad (3.15)$$

In tutti i casi la combinazione $([\]*)$ indica un numero variabile di parole intervenienti (zero compreso) che possono assumere il ruolo di soggetto verbale, ed è applicabile sono nel caso in cui si abbia annotazione sintattica che conservi il legame tra oggetto e verbo.

Per ognuno dei test di trasformazione sintattica è possibile considerare l'indice $I_{syn}^{(k)}$, con k a rappresentare topicalizzazione, ripresa anaforica, passivizzazione o relativizzazione, che valuti la proporzione delle occorrenze dell'espressione trasformata $n_{syn}^{(k)}$ rispetto a quelle della sua forma base. Per ciascuno degli indici, la formula di calcolo è la seguente:

$$I_{syn}^{(k)} = \frac{n_{syn}^{(k)}}{n_{bf} + n_{syn}^{(k)}} \quad (3.16)$$

3.3.2.4 Indice di modificazione sintagmatica

A seconda dei test di modificazione ammessi per ogni pattern, è possibile riassumere l'informazione di variazione in un unico indice globale I_{syn} nel seguente modo. Denotando con j ognuno dei test di modificabilità associato al pattern e con n_j il numero di occorrenze dell'espressione modificata secondo il test j -esimo, l'indice globale di modificabilità è dato da:

$$I_{syn} = \frac{\sum_j n_j}{n_{bf} + \sum_j n_j} \quad (3.17)$$

3.3.3 Variazioni paradigmatiche

Per ognuna delle espressioni della lista in input, lo strumento è in grado di compiere un test sulla variabilità paradigmatica, studiando la sostituibilità dei componenti dell'espressione con lemmi che siano in rapporto di sinonimia con essi. Tale test risulta l'unico di matrice semantica.

Al fine di verificare tale possibilità è necessario disporre di una risorsa esterna costituita da un tesoro di sinonimi. Si è scelto di considerare, nel presente lavoro, il tesoro italiano messo a disposizione dalla piattaforma GNU-OpenOffice⁵ grazie alla immediata disponibilità ai fini dell'utilizzo e alla facilità di utilizzo in quanto strutturato in formato testuale grezzo⁶. La risorsa annovera 26.823 lemmi i cui relativi sinonimi sono suddivisi per accezioni, come mostrato nell'esempio in Figura 3.1.

```
campo|8
(s.m.)|agro
(s.m.)|appezzamento|campagna|coltivazione|fondo|podere
(s.m.)|arena|pista|stadio
(s.m.)|accampamento|attendamento|bivacco
(s.m.)|area|regione|zona
(s.m.)|spazio
(s.m.)|sfondo|superficie
(s.m.)|ambito|argomento|branca|dominio|materia|ramo|sfera
```

Figura 3.1: Esempio di voce del tesoro OpenOffice per l'italiano. Nella prima riga è presente il lemma con il numero di accezioni per cui sono presenti sinonimi. Le righe successive indicano la categoria grammaticale del lemma e l'elenco di sinonimi afferenti ad una stessa accezione.

L'idea alla base del test di sostituibilità è quella di stabilire se un'espressione sia modificabile paradigmaticamente sulla base dell'attestazione nel corpus di espressioni ottenute sostituendo nell'originale ad uno dei suoi componenti ognuno dei relativi sinonimi. Maggiore sarà il numero di espressioni attestate mediante tale trasformazione, maggiore sarà la sostituibilità dell'espressione iniziale. Al contrario, nel caso di assenza (o di presenza trascurabile) nel corpus di espressioni ottenute dall'originale tramite sostituzione di uno dei componenti, l'espressione verrà considerata non modificabile paradigmaticamente. Anche in questo caso le query di ricerca di espressioni iniziali e varianti sostituite sono effettuate a livello di lemmi dei componenti. I componenti sottoposti a sostituzione sono unicamente le parole piene dell'espressione e le espressioni scaturite dalla trasformazione hanno, di volta

⁵http://linguistico.sourceforge.net/pages/thesaurus_italiano.html

⁶Non è esclusa, in ogni caso, la possibilità di un ricorso ad altre risorse, ove disponibili, al fine di disporre di una maggiore accuratezza in termini linguistici. Ad ogni modo, ai fini del presente studio, il tesoro considerato si è rivelato sufficientemente adatto alle analisi.

in volta, solo uno dei componenti sostituiti e mai entrambi contemporaneamente. I motivi di quest'ultima scelta verranno chiariti nel paragrafo 3.3.3.3.

3.3.3.1 Una versione naïf della sostituibilità empirica

Come primo approccio è possibile quantificare la variabilità paradigmatica sulla base della semplice attestazione delle espressioni sostituite. Ai fini esplicativi, si consideri l'esempio di *colonna sonora*, espressione riconducibile al pattern NA. Dal tesaurus a disposizione, i due componenti risultano possedere rispettivamente 24 e 8 sinonimi, come mostrato in Figura 3.2.

```
colonna|6
(s.f.)|pilastro|sostegno
(s.f.)|cariatide|cippo|obelisco|stele
(s.f.)|aiuto|appoggio|cardine|fondamento|perno|sostegno
(s.f.)|elenco|fila|serie
(s.f.)|carovana|coda|compagnia|drappello|fila|formazione|schiera
(s.f.)|banda|pista

sonoro|1
(agg.)|acustico|altisonante|enfatico|forte|risonante|roboante
      |rumoroso|squillante
```

Figura 3.2: Le voci *colonna* e *sonoro* del tesaurus italiano OpenOffice.

Lo strumento procede, quindi, alla sostituzione del primo lemma (*colonna*) con ognuno dei suoi sinonimi, ottenendo le espressioni lemmatizzate trasformate *pilastro sonoro*, *sostegno sonoro*, *cariatide sonoro*, ecc., di cui si ricerca il numero di occorrenze nel corpus. Esaurite le sostituzioni per il primo componente, si compiono analoghe trasformazioni per il secondo, ottenendo espressioni del tipo *colonna acustico*, *colonna altisonante*, e così via, calcolando anche per queste ultime il numero delle eventuali occorrenze nel corpus.

In generale, detto $syn_{1,i}$ l' i -esimo sinonimo di w_1 e $n_{syn_{1,i}}$ il numero di occorrenze dell'espressione sostituita della forma $syn_{1,i} w_2$ e, in maniera analoga $syn_{2,i}$ l' i -esimo sinonimo di w_2 e $n_{syn_{2,i}}$ il numero di attestazioni dell'espressione $w_1 syn_{2,i}$, il numero totale n_{sub}^{raw} di espressioni attestata in cui vi è sostituzione tramite sinonimo risulta:

$$n_{sub}^{raw} = \sum_i n_{syn_{1,i}} + \sum_i n_{syn_{2,i}} \quad (3.18)$$

Nel caso di pattern a tre componenti, di cui il centrale sia una parola vuota nella forma di preposizione o articolo, il calcolo è analogo a quello riportato sopra, a patto

di sostituire l'indice 2 con 3, dato che la forma dell'espressione iniziale risulta $w_1 w_2 w_3$.

Va precisato che lo strumento è in grado di riconoscere se il lemma centrale w_2 è una preposizione o un articolo soggetto a variazioni di cancellazione o inserzione fonetica, testando tutte varianti che potrebbero costituire dei lemmi autonomi con i possibili sinonimi presi in considerazione⁷ secondo lo schema in Tabella 3.1, riuscendo così a garantire la formazione di espressioni sostituite grammaticalmente corrette.

w_2	Varianti generate
a	a, ad
ad	a, ad
d'	d', da, di
da	d', da
di	d', di
e	e, ed
ed	e, ed
o	o, od
od	o, od
un'	un', un, una, uno
un	un', un, una, uno
una	un', un, una, uno
uno	un', un, una, uno

Tabella 3.1: Varianti considerate nella generazione di possibili espressioni sostituite per preposizioni e articoli soggetti a modifiche fonetiche.

Analogamente a quanto visto per i test di variabilità sintagmatica, è possibile costruire un indice che quantifichi la variabilità paradigmatica secondo la seguente formula:

$$I_{sub}^{raw} = \frac{n_{sub}^{raw}}{n_{bf} + n_{sub}^{raw}} \quad (3.19)$$

Anche in questo caso, l'indice ha variazione nell'intervallo compreso tra 0 e 1.

3.3.3.2 Quantificazione ragionata della sostituibilità empirica

L'approccio visto nel precedente paragrafo, benché utile per una prima individuazione di massima del grado di sostituibilità delle espressioni (cfr. Squillante 2014) presenta un forte limite metodologico. Se, infatti, la non attestazione di espressioni sostituite garantisce a livello empirico un'indicazione sull'impossibilità dell'entità di subire modifiche paradigmatiche, un'eventuale attestazione non è prova invece della

⁷Si pensi al caso di *specchio d'acqua*, che necessita della sostituzione di *d'* con *di* qualora si voglia testare la sequenza grammaticalmente corretta *specchio di liquido*, dove ad *acqua* è stato sostituito un suo sinonimo.

sua modificabilità. Per alcune espressioni sussiste la possibilità che, dopo la sostituzione di uno dei componenti, l'espressione trasformata risulti ancora largamente attestata nel corpus, ma con un significato diverso da quello iniziale.

Si pensi al caso di *braccio destro* e si ipotizzi che in un corpus l'espressione compaia unicamente in frasi che selezionino la sua lettura idiomatica, cioè quella di "aiutante" o "uomo di fiducia". Il primo componente dell'espressione ha, tra i suoi sinonimi, il lemma *ala* nell'accezione che raggruppa i termini che si riferiscono in generale alla parte di un edificio. Mediante la sostituzione di *braccio* con *ala*, si ottiene l'espressione *ala destra*, che, oltre alla lettura compositiva, designa il ruolo in campo di un giocatore di calcio. Come è evidente, i significati finali dell'espressione sostituita hanno poco in comune con il significato di partenza. Applicando in maniera cieca il test di sostituibilità nella forma delle equazioni 3.18 e 3.19, qualora *ala destra* fosse largamente attestato nel corpus, i risultati collocherebbero *braccio destro* fra le espressioni sostituibili.

Casi di questo tipo suggeriscono che i risultati di un test di sostituibilità affidabile non debbano dare risposta alla semplice domanda: "*Sostituendo uno dei componenti con un suo sinonimo, l'espressione risulta ancora attestata?*", bensì ad una più precisa: "*Sostituendo uno dei componenti con un suo sinonimo, l'espressione risulta ancora attestata nel suo significato originale?*".

Per rispondere a quest'ultima domanda e mantenere il carattere di procedura automatica dello strumento, si è reso necessario ricorrere ad un approccio basato sui metodi della *semantica distribuzionale*, che identificano e quantificano la similarità tra entità lessicali sulla base del loro distribuirsi in contesti simili o differenti.

Una panoramica esaustiva sulle motivazioni, lo sviluppo e la differenziazione di tali linee di ricerca e metodologie esula dagli scopi del presente lavoro. Tuttavia, per chiarezza, è qui utile ricordare le premesse teoriche ed epistemologiche di tali approcci, prima fra tutte la cosiddetta **ipotesi distribuzionale**, secondo cui *quanto più due entità lessicali occorrono in contesti linguistici simili, tanto più esse sono semanticamente simili* (Miller & Charles, 1991).

L'ipotesi si ricollega agli studi semantici dell'uso del linguaggio, sviluppati già negli anni cinquanta da Harris (1954, 1968) e Firth (1957), in cui emerge l'idea di studiare i comportamenti paradigmatici delle entità lessicali sulla base dei loro rapporti sintagmatici. Le parole con cui un'entità lessicale cooccorre creano una struttura d'informazione (vettore) che può essere confrontata con quella di un altro elemento del lessico. La similarità delle due strutture d'informazione fornisce il grado di interscambiabilità delle entità lessicali, come verrà esplicitato di qui a poco. Il significato è in correlazione diretta, quindi, ad uno *spazio di parole*, in cui il lessico è concepito come «uno spazio metrico i cui elementi - le parole - sono separati da distanze che dipendono dal loro grado di *similarità semantica*» (Lenci, 2009, p. 84).

In quest'ottica, anche caso dei fenomeni multiparola, determinare se due espressioni sono sinonime si traduce nel determinare quanto simili sono i contesti in cui esse occorrono.

guerra mondiale		conflitto mondiale		spedizione mondiale	
lemma	freq.	lemma	freq.	lemma	freq.
essere	17.287	essere	1.337	partecipare	4
avere	4.035	avere	425	avere	3
venire	3.533	venire	302	tuttavia	2
suo	3.415	suo	269	perdere	2
anno	2.568	anno	258	parte	2
militare	2.486	non	225	infortunio	2
tedesco	2.268	più	205	impedire	2
fine	2.201	anche	193	già	2
più	2.135	parte	189	fare	2
parte	2.073	fine	188	edizione	2
non	2.001	ultimo	146	arrivare	2
anche	1.957	scoppio	128	giungere	1
italiano	1.671	italiano	124	girone	1
aereo	1.477	italia	124	giocatore	1
scoppio	1.375	tedesco	120	ginocchia	1

Figura 3.3: Primi quindici lemmi più frequenti in cooccorrenza con le espressioni *guerra mondiale* (sinistra), *conflitto mondiale* (centro), *spedizione mondiale* (destra) nel corpus PAISÀ.

A tal fine, lo strumento opera la seguente procedura. Data un'espressione in input, viene generato un sottocorpus formato da tutte e sole le frasi in cui l'espressione è presente. Si calcolano, quindi, le frequenze di occorrenza di tutte le parole piene (aggettivi, sostantivi, verbi, avverbi) nel sottocorpus, ad esclusione dei lemmi componenti l'espressione. Tali frequenze, per costruzione, rappresentano nient'altro che le cooccorrenze delle parole piene con l'espressione multiparola in esame, dove lo *span* è assunto pari all'intera lunghezza della frase.

A titolo di esempio, nella prima tabella di Figura 3.3 sono riportati i quindici lemmi più frequenti che cooccorrono con l'espressione *guerra mondiale* nel corpus PAISÀ.

Il *vettore distribuzionale* che costituisce la struttura d'informazione relativa all'espressione in esame cui si è accennato sopra è, appunto, l'insieme delle frequenze di cooccorrenza di ciascun lemma con l'espressione. Per garantire l'agevolezza del carico computazionale nell'esecuzione della procedura automatica dello strumento, si è scelto di stabilire nel numero di 50 le componenti dei vettori presi in considerazione. Per ogni espressione, cioè, si conserva l'informazione di cooccorrenza dei soli 50 primi lemmi cooccorrenti più frequenti.

Come mostrato dalle altre due tabelle di Figura 3.3, lo strumento è in grado di calcolare i vettori distribuzionali per ciascuna delle espressioni generate per sostituzione di uno dei componenti con un sinonimo⁸, qualora l'espressione trasformata

⁸Per i componenti di *guerra mondiale* il tesoro suggerisce i seguenti sinonimi:

guerra belligeranza, campagna, conflitto, scontro, spedizione | antagonismo, confronto, contesa, contrasto, discordia, disputa, dissidio, ostilità, rivalità;

risultati attestata nel corpus.

Una volta ottenuti i vettori per tutte le espressioni trasformate, è necessario uniformare le loro componenti, in modo che ognuna di esse sia ordinatamente associata ad un unico lemma. A tal fine lo strumento genera una matrice distribuzionale $r \times c$, incrociando le informazioni vettoriali, in modo che in riga siano presenti i lemmi ed in colonna le frequenze di cooccorrenza con ciascuna espressione. Si noti che il numero di righe r è variabile, con la condizione $r \geq 50$, in quanto ognuno dei vettori generati dalle espressioni trasformate può introdurre un numero variabile di lemmi non presenti fra i 50 più frequenti per l'espressione originale, ma che ottengono comunque una propria riga all'interno della matrice. Un esempio di matrice generata dai tre vettori a 15 componenti di Figura 3.3 è mostrata in Tabella 3.2.

Si passa poi al confronto di similarità tra ognuno dei vettori delle espressioni trasformate (colonne $c \geq 2$ della matrice) e l'espressione originale ($c = 1$), al fine di ottenere un valore che rappresenti il grado di similarità semantica (sinonimia) tra l'espressione trasformata e quella originale.

Una delle misure più utilizzate in ambito distribuzionale per quantificare la similarità vettoriale, e che viene utilizzata in questo caso dallo strumento, è il **coseno di similitudine**, nota anche come **distanza di coseno**.

Dati due vettori \vec{A} e \vec{B} , di ugual numero n di componenti, la distanza di coseno è espressa dalla seguente formula:

$$\cos(\theta) = \frac{\vec{A} \cdot \vec{B}}{\|\vec{A}\| \|\vec{B}\|} = \frac{\sum_{k=1}^n A_k B_k}{\sqrt{\sum_{k=1}^n A_k^2} \sqrt{\sum_{k=1}^n B_k^2}} \quad (3.20)$$

dove A_k e B_k rappresentano le k -esime componenti dei vettori A e B , $\| \cdot \|$ è il simbolo di norma vettoriale e θ rappresenta l'angolo che separa i due vettori A e B nello spazio di parole. Poiché la distanza di coseno è espressa attraverso il calcolo del coseno dell'angolo θ , essa ha, in generale, la proprietà di essere una quantità limitata all'intervallo $[-1, 1]$. Nel nostro caso, però, poiché le componenti dei vettori sono espresse in termini di frequenze, e quindi di quantità positive, sarà possibile ottenere valori di coseno limitati all'intervallo che va da 0 a 1. Grazie a questa proprietà, sarà possibile interpretare il risultato della similarità di coseno come una *percentuale* di sinonimia distribuzionale tra i vettori presi in considerazione.

È importante notare che, poiché i vettori vengono normalizzati nel computo della loro distanza di similarità, la misura risulta indipendente dalla frequenza di occorrenza delle espressioni. Non conta, cioè, la “lunghezza” dei vettori rappresentanti le espressioni, bensì solo la loro orientazione nello spazio. La distanza di coseno tende, in altre parole, a verificare la similarità fra entità lessicali sulla base della similarità di proporzioni di distribuzione delle cooccorrenza dei lemmi in rapporto al numero

mondiale eccezionale, enorme, globale, intercontinentale, internazionale, universale.

In PAISÀ esistono attestazioni delle sole espressioni trasformate: *campagna/confitto/scontro/spedizione/confronto/contesa mondiale* e *guerra enorme/globale/intercontinentale/internazionale/universale*.

	guerra mondiale	conflitto mondiale	spedizione mondiale
aereo	1.477	0	0
anche	1.957	193	0
anno	2.568	258	0
arrivare	0	0	2
avere	4.035	425	3
edizione	0	0	2
essere	17.287	1.337	0
fare	0	0	2
fine	2.201	188	0
già	0	0	2
ginocchia	0	0	1
giocatore	0	0	1
girone	0	0	1
giungere	0	0	1
impedire	0	0	2
infortunio	0	0	2
italia	0	124	0
italiano	1.671	124	0
militare	2.486	0	0
non	2.001	225	0
parte	2.073	189	2
partecipare	0	0	4
perdere	0	0	2
più	2.135	205	0
scoppio	1.375	128	0
suo	3.415	269	0
tedesco	2.268	120	0
tuttavia	0	0	2
ultimo	0	146	0
venire	3.533	302	0

Tabella 3.2: Matrice dei vettori distribuzionali dei quindici lemmi più frequenti cooccorrenti con le espressioni *guerra mondiale*, *conflitto mondiale* e *spedizione mondiale* nel corpus PAISÀ. Come si vede il numero totale di righe (30) è maggiore del numero di componenti iniziali (15).

totale di cooccorrenze. Grazie a questa caratteristica non risulta problematico il raffronto tra vettori che rappresentino espressioni con occorrenze molto diverse tra loro⁹.

Un calcolo della distanza di coseno tra i vettori rappresentanti le espressioni di Figura 3.3, estesa a 50 componenti, fornisce una percentuale di sinonimia del 97% tra *guerra mondiale* e *conflitto mondiale*, mentre solo dell'11% tra *guerra mondiale* e *spedizione mondiale*.

Una volta computate le distanze di coseno, è possibile sfruttare la percentuale di similarità tra le espressioni trasformate e l'espressione originale come un fattore peso per il calcolo delle occorrenze di attestazione delle espressioni trasformate. L'idea è quella di moltiplicare il numero di occorrenze dell'espressione sostituita nel corpus con la percentuale di similarità rispetto all'espressione che la ha generata. In questo modo, le occorrenze di espressioni ad alta sinonimia verranno normalizzate ad un numero pressoché simile; quelle di espressioni a bassissimo grado di sinonimia con l'originale risulteranno in un numero di occorrenze normalizzato quasi nullo; le espressioni con percentuali di sinonimia intermedie vedranno il numero delle proprie attestazioni normalizzarsi in proporzione al loro grado di sinonimia con l'originale.

In definitiva, lo strumento sarà in grado di calcolare un numero di attestazioni delle espressioni sostituite opportunamente modificato rispetto a quello dell'equazione 3.18, dato dalla seguente formula:

$$n_{sub} = \sum_i \cos \theta_{1,i} n_{syn_{1,i}} + \sum_i \cos \theta_{2,i} n_{syn_{2,i}} \quad (3.21)$$

dove i fattori coseni indicano la distanza di coseno tra l'espressione sostituita tramite sostituzione dell'i-esimo sinonimo del primo o del secondo componente e l'espressione originale.

L'indice di sostituibilità I_{sub} viene quindi espresso secondo la formula:

$$I_{sub} = \frac{n_{sub}}{n_{bf} + n_{sub}} \quad (3.22)$$

Anche se non ne verrà fatto uso attivo nelle successive analisi, per completezza, in aggiunta a tale indice, può risultare interessante avere anche una stima della dispersione della sostituibilità dei componenti, vale a dire un'informazione in grado di stabilire se il numero totale di occorrenze delle espressioni con sostituzione di uno stesso componente derivi da occorrenze distribuite in maniera più o meno uniforme sui diversi sinonimi del componente o concentrate principalmente solo su uno o pochi di essi.

⁹È possibile, tuttavia, che il basso numero di occorrenze di un'espressione limiti la variabilità dei lemmi in cooccorrenza, introducendo un bias nel confronto con un'altra espressione che, ampiamente attestata, può vantare una distribuzione di lemmi in cooccorrenza più varia ma soprattutto più robusta. Va ricordato, in ogni caso, che espressioni trasformate che abbiano un esiguo numero di occorrenze rispetto all'espressione originale, hanno un peso minimo nel calcolo empirico della sostituibilità, come specificato più avanti dall'equazione 3.21.

Si definisce quindi la dispersione D_k di sostituibilità di un componente w_k come:

$$D_k = \frac{s_k}{\sigma_k} \quad (3.23)$$

con $k = 1, 2$, dove s_k è il numero dei sinonimi di w_k per i quali l'espressione interessata dalla sostituzione ha avuto un numero di occorrenze diverse da zero e σ_k la deviazione standard della distribuzione delle occorrenze tra i sinonimi attestati, di cui va però precisato il significato.

Se $n_{syn_1} = \sum_i \cos \theta_{1,i} n_{syn_{1,i}}$ è la somma normalizzata delle occorrenze delle espressioni che vedono la sostituzione del primo componente con i suoi sinonimi e s_1 il numero di sinonimi che, se sostituiti al primo componente, forniscono espressioni attestate nel corpus, si consideri il valore medio di occorrenze

$$m_1 = \frac{n_{syn_1}}{s_1} \quad (3.24)$$

che rappresenta il numero occorrenze nel caso in cui esse siano equamente distribuite tra tutti i sinonimi attestati.

Nel nostro caso σ_1 è calcolata nella maniera canonica:

$$\sigma_1 = \sqrt{\frac{\sum_i (n_{syn_{1,i}} - m_1)^2}{s_1}} \quad (3.25)$$

ed è tanto più bassa quanto più le occorrenze delle espressioni relative alla sostituzione di ogni sinonimo si avvicinano al caso in cui esse siano equamente distribuite tra tutti i sinonimi attestati.

Dalla equazione (3.23) è possibile vedere che il valore di D_k aumenta all'aumentare del numero di sinonimi s_k attestati dalle espressioni nel corpus. Allo stesso tempo esso diminuisce se σ_k è un valore alto, il che equivale a dire che le occorrenze delle espressioni con sostituzione sono divise solo tra uno o pochi dei sinonimi attestati.

3.3.3.3 Limiti del test di sostituibilità empirica

Per quanto la rinormalizzazione nel calcolo delle occorrenze attestate dell'espressione sostituita renda più robusti i valori dell'indice I_{sub} di sostituibilità empirica, sussistono dei limiti metodologici in tale approccio.

Un primo problema sussiste per le espressioni multiparola che possiedono espressioni sinonime ottenibili per sostituzione di uno dei componenti, ma la sostituzione non avviene tramite sinonimo. Un tipico esempio di questa situazione è il caso di *campo di concentramento*, sinonimo dell'espressione *campo di sterminio*, benché *concentramento* e *sterminio* non stiano in un rapporto di sinonimia, bensì di consequenzialità. Come evidente da quanto esposto nell'algoritmo procedurale, lo strumento non è quindi in grado di generare espressioni trasformate con lemmi che non siano sinonimi dei componenti iniziali. Ciononostante è possibile superare empiricamente

tale ostacolo attraverso l'inclusione nell'insieme dei lemmi da sostituire oltre che dei sinonimi di ciascun componente, anche i lemmi che più frequentemente cooccorrono con l'altro componente. Nel caso di *campo di concentramento*, è possibile infatti effettuare nel corpus una ricerca delle parole maggiormente cooccorrenti al lemma *campo* e ordinarle per frequenza. Fissato un numero massimo di lemmi da considerare (*n-best list*) è possibile ricercare nel corpus le attestazioni dell'espressione trasformata attraverso l'utilizzo di ognuno dei lemmi frequenti, se non già inclusi nell'insieme dei sinonimi. Il grado di sinonimia risultante dal calcolo della distanza di coseno filtrerà, quindi, le occorrenze "utili" perché relative ad un'espressione rivelatasi sinonima dell'originale, da quelle da scartare. In relazione all'esempio citato, i sette sostantivi più frequenti a completamento dell'espressione *campo di* nel corpus PAISÀ risultano: *concentramento* (2.360), *battaglia* (2.146), *gioco* (786), *sterminio* (479), *prigionia* (440), *ricerca* (379), *lavoro* (370). Tra questi, quelli generanti le espressioni trasformate più vicine distribuzionalmente a *campo di concentramento* sono *campo di sterminio* (94%), *campo di prigionia* (89%) e *campo di lavoro* (87%), mentre le restanti espressioni trasformate rimangono sotto il 40% di similarità. Sviluppi futuri dello strumento potranno, quindi, tenere in considerazione tale opzione.

Infine un secondo limite è relativo al fatto che il programma genera espressioni a partire dalla sostituzione di un solo componente alla volta, tralasciando la possibilità di sostituzione contemporanea (si pensi al caso di *presidente della repubblica/capo dello stato*). Questa scelta è stata dettata primariamente dall'agevolezza del carico computazionale. Si supponga, ad esempio, che il primo componente abbia n sinonimi e il secondo componente ne abbia m . Mentre con l'attuale struttura il numero di espressioni da testare risulta nel numero di $n + m$, nel caso di sostituzione contemporanea tale cifra salirebbe a $n \cdot m$. Alcuni test effettuati hanno inoltre mostrato come la sostituzione contemporanea dei componenti con possibili sinonimi conduca ad un'ingente dispersione del significato originale dell'espressione, come nel caso di *via d'uscita*, per cui si arriva a generare *inizio di pubblicazione* o *itinerario d'apertura*. Nonostante in questi casi il calcolo della similarità distribuzionale potrebbe ridimensionare le occorrenze delle espressioni non sinonime, gli ingenti tempi e risorse di calcolo per tenere in conto tale possibilità sfavoriscono un tale approccio.

3.3.4 Variazioni flessive

Per ognuna delle espressioni della lista iniziale che, ricordiamo, sono inserite in input con i componenti lemmatizzati, lo strumento è in grado di valutare se esista nel corpus una concentrazione di occorrenze in una particolare forma flessa dei componenti. Il test è utile ad individuare i casi in cui c'è, appunto, cristallizzazione flessiva dell'espressione: si pensi al caso di *fai da te* del pattern VPN o di *acque bianche* per il pattern NA.

Data l'espressione, lo strumento conta le frequenze di occorrenza attestate nel corpus per ciascuna sua combinazione di forme flesse, ordinandole in maniera de-

crescente. Detto n_{prev} il numero di occorrenze della forma flessa più frequente (la prima della lista), l'indice di variabilità flessiva I_{infl} è calcolato nel seguente modo:

$$I_{infl} = \frac{n_{bf} - n_{prev}}{n_{bf}} \quad (3.26)$$

Anche in questo caso, l'indice è vincolato a produrre valori compresi nell'intervallo da 0 a 1.

Come si vede, l'indice valuta la proporzione del numero totale di occorrenze delle forme flesse esclusa quella prevalente rispetto al numero totale di occorrenze dell'espressione. Di conseguenza, tanto più le occorrenze dell'espressione sono concentrate in un'unica forma, tanto minore sarà il valore di I_{infl} , attestando, appunto, una bassa tendenza dell'espressione ad essere flessa.

Per completezza, il programma fornisce anche l'informazione sul numero di occorrenze della seconda forma flessa più frequente, in modo da poter valutare il distacco in termini di frequenza tra prima e seconda forma prevalente ed, eventualmente, quanto le occorrenze delle prime due forme flesse coprano delle occorrenze totali.

3.4 Una scelta metodologica sulla valutazione della categorizzazione

Come si vedrà nel prossimo capitolo, l'analisi dei risultati degli indici di variabilità forniti dallo strumento mostrerà la possibilità di categorizzare il continuum interno ai fenomeni multiparola, oltre che quella tra questi ultimi e le espressioni libere, in polarizzazioni che individuano categorie già note alla teoria linguistica. Queste ultime, tuttavia, assumeranno una caratterizzazione fissata da criteri empirici, che permetteranno di fornirne una definizione più precisa.

In quest'ottica la questione della validazione dei risultati e della bontà della categorizzazione diventa un terreno delicato. Da un lato, infatti, si preferirebbe evitare l'imposizione di categorie *a priori*, dandone delle definizioni prototipiche su base intuitiva e mirando poi ad ottenere giudizi di annotazione da parte di annotatori, grazie ai quali si possa effettuare un confronto con le polarizzazioni individuate dallo strumento. In questo caso, oltre all'aleatorietà della scelta delle caratteristiche pertinenti alla definizione di una certa classe di espressioni e alla problematicità di collocazione delle entità che non mostrano palese riconducibilità alla classe, esiste anche un problema di estensione delle definizioni oltre il confine di gruppi ristretti di categorie grammaticali o pattern.

D'altro canto, senza categorie stabilite *a priori* il lavoro di evidenziazione di classi di espressioni differenziate in base ai risultati forniti dal programma e della loro definizione è demandato a chi analizza i dati e alla sua capacità di identificare con obiettività eventuali omogeneità di caratteristiche delle espressioni che mostrano comportamenti empirici comuni.

Un possibile rimedio sarebbe il ricorso alle risorse lessicografiche per l'italiano, quali ad esempio l'insieme di polirematiche del GRADIT (De Mauro, 1999-2007) e i dizionari combinatori di Urzì (2009), Lo Cascio (2011) e Tiberii (2012) al fine di poter avere una base ufficiale su cui operare il discernimento categoriale tra polirematiche (per cui il GRADIT rappresenta oggi ancora la più grande risorsa lessicografica per l'italiano) e collocazioni (individuate come preferenze combinatorie degli altri tre dizionari).

Poiché, tuttavia, i confini teorici di identificazione di tali classi sono confusi e diversi a seconda dei lavori che, spesso, hanno scopi applicativi e orientati all'utente generico più che al linguista, le stesse risorse lessicografiche appaiono non affidabili a garantire un riferimento per il discernimento in classi delle espressioni.

A titolo di esempio, delle 500 espressioni più frequenti afferenti al pattern NA del corpus PAISÀ, 386 risultano attestate in almeno uno dei quattro dizionari sopra elencati¹⁰. Di queste, 182 risultano attestate nel GRADIT, e 135 sono voci di tutti e tre i dizionari combinatori. I due insiemi, tuttavia, sono sovrapposti, in quanto ben 172 espressioni sono contemporaneamente attestate sia nel GRADIT che in almeno uno dei tre dizionari di collocazioni, mostrando quanto i dizionari non abbiano una chiara strategia di distinzione tra espressioni unitarie o cristallizzate che necessitano della combinazione dei costituenti per veicolare il proprio significato, da quelle per cui i componenti esibiscono solo preferenzialità di combinazione. Del resto, se il GRADIT mira ad annoverare di norma le unità polirematiche¹¹ i dizionari combinatori identificano legami lessicali per più generiche espressioni fraseologiche non fisse ma riconoscibili (Tiberii, 2012), per sintagmi caratterizzati da un grado più o meno accentuato di coesione interna (Urzì, 2009) o, più in generale, per famiglie di parole (Lo Cascio, 2011).

In virtù delle considerazioni esposte, quindi, si è scelto di optare per la seconda fra le opzioni viste in precedenza, ovvero per un'analisi **qualitativa** che, a valle dei risultati forniti dallo strumento, provi ad individuare sulla base dell'omogeneità delle polarizzazioni le caratteristiche empiriche delineanti le classi categoriali.

¹⁰Si ricordi che tra le espressioni più frequenti di un pattern è naturale che siano comprese anche espressioni libere.

¹¹Secondo la definizione di De Mauro, le espressioni incluse nella risorsa devono soddisfare almeno uno dei seguenti requisiti:

- l'esistenza di uno specifico sovrappiù semantico, vale a dire la non ricostruibilità del loro significato in base alla semplice somma dei significati dei singoli componenti monorematici [...];
- la più o meno forte cristallizzazione lessicale e sintattica, ovvero il fatto che la polirematica, in quanto considerata come un unico elemento lessicale, tende a non ammettere variazioni lessicali e strutturali interne senza che si perda il sovrappiù semantico di cui è portatrice [...];
- la presenza significativa in linguaggi tecnico-specialistici [...].

(De Mauro, 2005).

Analisi sul linguaggio generale dell'italiano

4.1 Il corpus PAISÀ

Una volta stabilita la metodologia d'indagine della variabilità morfologica e lessico-sintattica delle espressioni multiparola e costruito lo strumento computazionale in grado di automatizzare la produzione dei dati risultanti dalle query, si è scelto di utilizzare il corpus PAISÀ¹ (Piattaforma per l'Apprendimento dell'Italiano Su corpora Annotati, Lyding *et al.* 2014) come risorsa di riferimento per l'analisi del linguaggio generale italiano.

PAISÀ rappresenta uno dei numerosi esempi contemporanei del consolidato utilizzo del World Wide Web come fonte per la creazione di risorse linguistiche². Esso, infatti, è costituito da una vasta raccolta di testi in lingua italiana estratti dal web e contrassegnati da alcune tipologie di licenze che permettono il loro libero utilizzo e rielaborazione³. La provenienza dei documenti è varia e distribuita secondo i dati in Tabella 4.1.

Affinché la risorsa fosse rappresentativa di un italiano comune e non orientato ad un qualche settore tecnico-specialistico, il reperimento dei documenti è stato basato su una *seed list* di lessemi comuni, ovvero una lista di combinazioni di parole che avrebbero dovuto essere presenti nei testi selezionati, ottenuta, come precisato dagli autori⁴, incrociando lemmi tratti dal Vocabolario di Base della lingua italiana (De Mauro, 1980), per un totale di 50.000 combinazioni. I documenti così individuati hanno subito quindi un processo di pulizia per eliminare eventuali sezioni di contenuto spurio e ritenuto non interessante a livello linguistico quali pubblicità, riferimenti e link, indici, ecc. Come esposto nel seguito, tuttavia, la natura automatica di tale processo non ha evitato il permanere, nel corpus, di segmenti ripetuti

¹www.corpusitaliano.it

²In relazione a questi usi si parla, generalmente, di *web as a corpus* (Kilgarriff & Grefenstette, 2001).

³Nello specifico, i testi sono sotto licenza Creative Commons (creativecommons.org). Per approfondire si rimanda a Lyding *et al.* (2014, p. 37).

⁴<http://www.corpusitaliano.it/it/contents/construction.html>

Fonte	N. documenti
Wikipedia	263.300
Wikibooks	2.380
Wikinews	1.680
Wikiversity	740
Wikisource	410
Wikivoyage	390
guide.supereva.it	19.000
italy.indymedia.org	10.000
tvblog.it	9.088
motorblog.it	3.300
ecowebnews	3.220
webmasterpoint.org	3.138
<i>altri</i>	71.250

Tabella 4.1: Ripartizione per provenienza dei documenti del corpus PAISÀ. La categoria *altri* racchiude documenti estratti da più di 1.000 siti diversi (Lyding *et al.*, 2014).

che hanno comportato l'esclusione manuale di alcune espressioni dall'analisi condotta in questa sede. Infine il corpus è stato sottoposto ad un'annotazione su tre livelli che attribuisse ad ogni token il proprio lemma, una categoria grammaticale⁵, e un inquadramento sintattico di tipo *dependency*⁶.

La scelta di tale corpus quale riferimento per le analisi compiute nel presente lavoro è giustificata da molteplici fattori. In primis, il corpus costituisce una risorsa di grandi dimensioni per l'italiano. Tale caratteristica tende a garantire, quindi, la vasta base empirica di cui si è discusso in precedenza, favorendo l'attendibilità dei risultati, specie considerando che le analisi sono svolte sulle espressioni più frequenti per ogni pattern. In secondo luogo PAISÀ campiona un italiano *contemporaneo*, in quanto tutti i documenti che costituiscono il corpus sono stati raccolti tra il settembre e l'ottobre del 2010 (Lyding *et al.*, 2014, p. 37). Le indagini beneficiano quindi di una base empirica attuale che possa contemplare espressioni (e relativi comportamenti linguistici) tra le più attuali. PAISÀ, inoltre, è liberamente disponibile e ripubblicabile in virtù dei criteri adottati nel collezionamento dei testi. La piena accessibilità favorisce dunque un suo utilizzo da parte dei linguisti nonché permette un controllo di riproducibilità su dati e analisi condotte su di esso.

⁵L'attribuzione di lemmi e *part of speech* è stata effettuata grazie al tagger descritto in Del-Orletta (2009), che combina sei algoritmi decisionali diversi al fine di ottenere le migliori performance di accuratezza (96,34% su PAISÀ, Lyding *et al.* 2014). Il *tagset*, ovvero l'insieme delle etichette grammaticali disponibili per l'annotazione è l'ISST-TANL, disponibile all'indirizzo: <http://www.corpusitaliano.it/static/documents/POS-ISST-TANL-tagset-web.pdf>, nonché in Appendice A.

⁶L'annotazione è stata effettuata grazie al DeSR Parser (Attardi *et al.*, 2009), secondo il *tagset* ISST-TANL disponibile all'indirizzo: <http://poesix1.ilc.cnr.it/ISST-TANL-DEPtagset-web.pdf>.

Infine, una delle motivazioni più importanti che ha governato la sua scelta in questa sede è la quantità di annotazione linguistica presente. I suoi tre livelli di annotazione, infatti, rendono massima l'accuratezza delle query effettuate dallo strumento computazionale nel reperire le espressioni, siano esse in forma base o modificate, come sarà evidente dalle analisi sui pattern esposte nei successivi paragrafi.

4.2 Analisi sui pattern nominali

È noto come alcune sequenze grammaticali siano quantitativamente più produttive nella creazione di espressioni di interesse fraseologico. In particolare, esiste una serie di pattern già individuati (Voghera, 2004; Masini, 2009) come i più comuni generatori di polirematiche e collocazioni che abbiano un comportamento analogo (o sostituibile) a quello dei sostantivi e che in questa sede verranno definiti *pattern nominali*. Può essere utile, a fini discorsivi, definire il pattern generatore di una certa espressione fraseologica *sequenza grammaticale in entrata* ed indicare la categoria funzionale dell'espressione così formata come *categoria grammaticale in uscita* per distinguere i livelli grammaticali di formazione e realizzazione dell'espressione.

La Tabella 4.2 riassume i pattern nominali analizzati nel presente lavoro, mostrando esempi di espressioni che essi sono in grado di formare.

Pattern	Esempio	Test I_{syn}
NA	<i>pollice verde</i>	int, ord
AN	<i>libero arbitrio</i>	int, ord
NPN	<i>luna di miele</i>	int
NPdN	<i>vigile del fuoco</i>	int
NPV _{inf}	<i>gomma da masticare</i>	int
NN	<i>scheda madre</i>	int
NCN	<i>botta e risposta</i>	int, ord
VCV	<i>gratta e vinci</i>	int, ord

Tabella 4.2: Lista dei pattern grammaticali di generazione di espressioni nominali analizzate nel presente lavoro. Per ogni pattern è mostrato un esempio di espressione generata e i test di variazione sintagmatica utilizzati nella generazione dei dati (int = interrompibilità; ord = ordine inverso).

L'ultima colonna della tabella indica anche quali proprietà di modificabilità sintagmatica vengono testate nella generazione dei valori per l'indice I_{syn} . A differenza del test di sostituibilità, infatti, la modificabilità sintagmatica può essere generata attraverso diverse trasformazioni, come precisato nel paragrafo 3.3.2. Per le espressioni nominali non è possibile effettuare i test di trasformazione sintattica verbale (ripresa anaforica, topicalizzazione, relativizzazione, passivizzazione). Inoltre non è sempre ragionevole verificare se le espressioni appartenenti a determinati pattern ammettano inversione dell'ordine dei costituenti. Per i pattern NPN, NPdN, NPV_{inf}

e NN l'inversione delle parole piene causa una modifica del significato originale (*casa di cura* → *cura di casa*; *vigile del fuoco* → *fuoco del vigile*) o la generazione di espressioni prive di significato (*associazione a delinquere* → *delinquere ad associazione*; *gas serra* → *serra gas*) e per questo in relazione ad essi il solo test di modifica sintagmatica operato è il test di interrompibilità.

Va precisato inoltre che, nelle analisi esposte di seguito, l'individuazione delle 500 espressioni più frequenti per ogni pattern costituenti la lista in input per lo strumento, è generata automaticamente. A valle della produzione dei dati, tuttavia, si è resa necessaria l'eliminazione di alcune espressioni a causa di tre principali e diverse ragioni.

In primo luogo esistono, nel corpus, casi di errata lemmatizzazione (ad es. *quartier* in luogo di *quartiere* nell'espressione *quartier generale*) o di errata attribuzione della corretta categoria grammaticale (ad es. *fine/A stagione/N*).

In secondo luogo esistono espressioni duplicate un ingente numero di volte nel corpus. Ciò è dovuto alla presenza delle espressioni all'interno di unità di testo che fanno parte, generalmente, di intestazioni, indici o sezioni di pagine web costituenti la parte strutturale del sito e non il contenuto della pagina. Tali espressioni, in virtù del loro presentarsi in maniera identica per un numero elevato di volte, ottengono tipicamente valori di modificabilità quasi nulla, ma non possono per questo essere rappresentative dell'insieme delle espressioni cristallizzate. In ragione di ciò, si è scelto di escludere dall'analisi espressioni che presentassero più del 50% delle occorrenze ripetute.

La terza motivazione di esclusione di espressioni dalle analisi è l'appartenenza dell'espressione, per più del 50% delle occorrenze, ad un nome proprio. È questo il caso, ad esempio, di sporadiche espressioni del tipo *posto al sole*, *distretto di polizia* o *mamma per amica* facenti parte di titoli di note serie televisive italiane.

Per le espressioni formate da lessico straniero, infine, si è scelto di lasciare nelle analisi le combinazioni pienamente integrate nella lingua solo se i componenti risultano etichettati grammaticalmente in maniera corretta (come nel caso di *social/A network/N*), mentre vengono esclusi i casi di errata categorizzazione (*hard/N disk/N*).

L'analisi dei dati empirici a valle dell'esecuzione dei test ha, inoltre, portato alla decisione di escludere l'indice di modificabilità flessiva quale asse di variabilità utile alla categorizzazione delle espressioni nominali. Come già mostrato in Squillante (2014) il blocco morfologico non sembra aggiungere maggiore informazione sullo status fraseologico dell'espressione rispetto a quanto già fornito dalle cristallizzazioni sul piano sintagmatico e paradigmatico. Tale conclusione può essere giustificata considerando che le espressioni nominali ammettono, come varianti flessive, le due sole forme di singolare e plurale e possono, perciò, ricadere unicamente nelle seguenti tre situazioni.

Se l'indice di variabilità flessiva presenta valori medi o alti (per semplicità si consideri $I_{infl} > 0,1$), ciò implica che l'espressione risulta attestata in entrambe le forme: tale situazione è però tipica sia di espressioni libere (ad es. *quanto giallo/-*

quantì gialli) che di espressioni tipicamente identificate come polirematiche (*cartone animato/cartoni animati*); queste ultime, infatti, configurandosi come entità sostituibili funzionalmente ad un unico sostantivo, acquisiscono la necessità di ammettere singolare e plurale.

Esistono, in secondo luogo, espressioni che mostrano assenza empirica di forma plurale (o una preponderanza considerevole della forma singolare). In questi casi, la flessione (plurale) bloccata non risulta indicatore di particolari legami fraseologici tra i componenti in quanto è possibile che la testa stessa dell'espressione non ammetta plurale per ragioni semantiche. È questo il caso, ad esempio, dei sostantivi astratti (come *bontà, felicità, coraggio*, ecc.) generalmente attestati nella sola forma singolare, ma non per questo generanti polirematiche o collocazioni se abbinati ad altro materiale linguistico. In modo analogo, un'entità fortemente terminologica come *anidride carbonica* non ammetterà plurale per le stesse ragioni semantiche piuttosto che per la sua natura di polirematica.

Il terzo caso, infine, è relativo a quelle espressioni che presentano un blocco empirico sulla forma singolare, essendo attestate in massima parte o unicamente al plurale ($I_{infl} < 0,1$). Anche in questa situazione esistono delle motivazioni semantiche che possono spiegare la rigidità flessiva, in particolare per il pattern NCN, dove sono comuni espressioni riferite a classi di individui (*cattolici e protestanti, artisti e intellettuali, piante ed animali*, ecc.). Tuttavia, per i pattern che generano espressioni che funzionalmente sono sostituibili a sostantivi, tralasciando i casi di parole italiane che già da sole risultano grammaticalmente dei *pluralia tantum*⁷, l'attestazione del solo plurale può essere indice di un certo grado di fissità fraseologica (si pensi a *giochi olimpici, diritti umani, acque bianche*, ecc.). Tuttavia tale blocco risulta sempre associato, a livello empirico, ad almeno un altro blocco variazionale. Per il pattern NA, ad esempio, il 100% delle espressioni la cui forma preferita è il plurale e che hanno una variabilità flessiva inferiore all'1%, hanno anche valori dell'indice di modificabilità sintagmatica inferiori al 4%, come mostrato in Tabella 4.3. In questi casi, data l'alta correlazione tra blocco flessivo e sintagmatico, è la modificabilità paradigmatica a differenziare la categoria dell'espressione, come mostrato nei paragrafi seguenti.

⁷Ad es. *forbici, occhiali, pantaloni, broccoli*, ecc. In ogni caso tali parole ammettono comunque la forma singolare in variazioni sia semantiche che diafasiche/diastratiche (in PAISÀ sono riscontrabili le seguenti occorrenze: *forbice* 306, *forbici* 410; *occhiale* 58; *occhiali* 2105; *pantalone* 169, *pantaloni* 1538; *broccolo* 11, *broccoli* 90).

Espr.	Forma preval.	Freq.	I_{syn}	I_{sub}	I_{infl}
problema economico	problemi economici	824	0,036257	0,513628	0,040049
diritto umano	diritti umani	3.081	0,002913	0,437834	0,017851
reperto archeologico	reperti archeologici	923	0,003240	0,271848	0,046587
bene culturale	beni culturali	735	0,016064	0,242217	0,048980
dato personale	dati personali	1.089	0,004570	0,215689	0,008264
elezione politico	elezioni politiche	2.647	0,003764	0,213747	0,003778
forza armato	forze armate	3.882	0,007922	0,205033	0,086296
ente locale	enti locali	1.220	0,008936	0,203349	0,100000
difficoltà economico	difficoltà economiche	1.122	0,002667	0,193274	0,055258
risorsa naturale	risorse naturali	890	0,005587	0,188183	0,048315
condizione ambientale	condizioni ambientali	774	0,015267	0,127113	0,027132
caratteristica tecnico	caratteristiche tecniche	2.514	0,001192	0,116662	0,001591
paese occidentale	paesi occidentali	766	0,026684	0,104162	0,039164
gioco olimpico	giochi olimpici	754	0,005277	0,087937	0,001326
elezione regionale	elezioni regionali	813	0,015738	0,075169	0,011070
regione temperato	regioni temperate	1.176	0,005076	0,069258	0,007653
mese estivo	mesi estivi	736	0,002710	0,068007	0,009511
prodotto agricolo	prodotti agricoli	744	0,005348	0,064605	0,063172
casa automobilistico	case automobilistiche	1.230	0	0,059178	0,060976
truppa tedesco	truppe tedesche	835	0,014168	0,058956	0,005988
invasione barbarico	invasioni barbariche	779	0,002561	0,055068	0,035944
condizione climatico	condizioni climatiche	797	0,007472	0,054315	0,040151
effetto speciale	effetti speciali	1.701	0,009319	0,041788	0,049383
arma nucleare	armi nucleari	982	0,008081	0,033571	0,091650
scienza naturale	scienze naturali	742	0,005362	0,024512	0,063342
elezione presidenziale	elezioni presidenziali	1.540	0,007092	0,017151	0,054545
truppa francese	truppe francesi	855	0,032805	0,011744	0,005848
consorzio agrario	consorzi agrari	1.238	0	0,011737	0,081583
aereo militare	aerei militari	2.808	0,000356	0,008086	0,029558
giorno scorso	giorni scorsi	1.143	0,041107	0,002567	0,003500
elezione europeo	elezioni europee	846	0,007042	0,000955	0,004728
scavo archeologico	scavi archeologici	980	0,002037	0,000156	0,077551

Tabella 4.3: Espressioni relative al pattern NA che presentano valori di variazione flessiva inferiori all'1%.

Un caso a parte è costituito, infine, dal pattern VCV che, nonostante la composizione delle categorie in entrata, è in grado di generare espressioni nominali (*gratta e vinci, tira e molla*). In questi casi l'assenza di flessione (che è qui intesa nelle varianti di persona, tempo, aspetto e modo) è imprescindibile requisito (come anche ipotizzato in De Mauro & Voghera, 1996) per la transcategorizzazione verso la categoria nominale. I pochi esempi di questo tipo ottenuti nelle analisi (cfr. Tabella 4.32) hanno mostrato come la cristallizzazione flessiva implichi anche la totale cristallizzazione sia sintagmatica che paradigmatica.

Alla luce di quanto detto la variabilità flessiva diviene, quindi, a tutti gli effetti un criterio secondario, e per questo trascurabile ai fini della categorizzazione delle espressioni.

Nei paragrafi seguenti sono esposte le analisi per ognuna delle sequenze grammaticali in entrata esposte in Tabella 4.2. Per comodità sono qui incluse la Figure 4.1 e 4.2, che mostrano la distribuzione delle 500 espressioni in input per ogni pattern nominale (una volta che quelle da scartare siano state filtrate) nello spazio generato dagli assi relativi agli indici I_{syn} e I_{sub} .

Per tutti i pattern esiste un nucleo di espressioni concentrato nell'angolo basso a sinistra, rappresentante l'area di cristallizzazione sintagmatica e paradigmatica.

Per alcuni pattern, poi, vi è una concentrazione delle espressioni sull'asse della non modificabilità sintagmatica. Tale caratteristica è particolarmente marcata per i pattern NA e NPN (e in misura minore per AN), che quindi mostrano una tendenza propria della combinazione sintattica ad opporre maggiore resistenza rispetto a interruzioni o inversione dei costituenti. Il pattern NPdN, invece, vede una concentrazione delle espressioni su una fascia di medio-bassa modificabilità sintagmatica, mostrando una maggiore flessibilità giustificabile dalla presenza dell'articolo determinativo, come precisato nel paragrafo 4.2.4.

L'alta concentrazione delle espressioni NPV_{inf} su valori di bassa modificabilità sintagmatica ma medio-alta variabilità paradigmatica è giustificata dal fatto che, come si vedrà nel paragrafo 4.2.5, la quasi totalità delle entità non cristallizzate è costituita da frammenti della sequenza PNPV_{inf} (es. *[in] grado di V_{inf}*), che quindi ammette sostituzione del verbo, ma richiede la fissità della costruzione.

Il pattern NN, non costituendo una sequenza integrata nella sintassi italiana, ma rappresentativa, invece, di entità assimilabili ai composti (cfr. paragrafo 4.2.6) mostra poche espressioni al di fuori dell'angolo di totale cristallizzazione, con una maggiore distribuzione sull'asse della sostituibilità rispetto a quello dell'interruzione.

I pattern in coordinazione (NCN e VCV), invece, sono i soli a distribuirsi con maggiore copertura sull'intero piano variazionale, in conseguenza della maggiore indipendenza sintattica che la congiunzione concede ai componenti pieni.

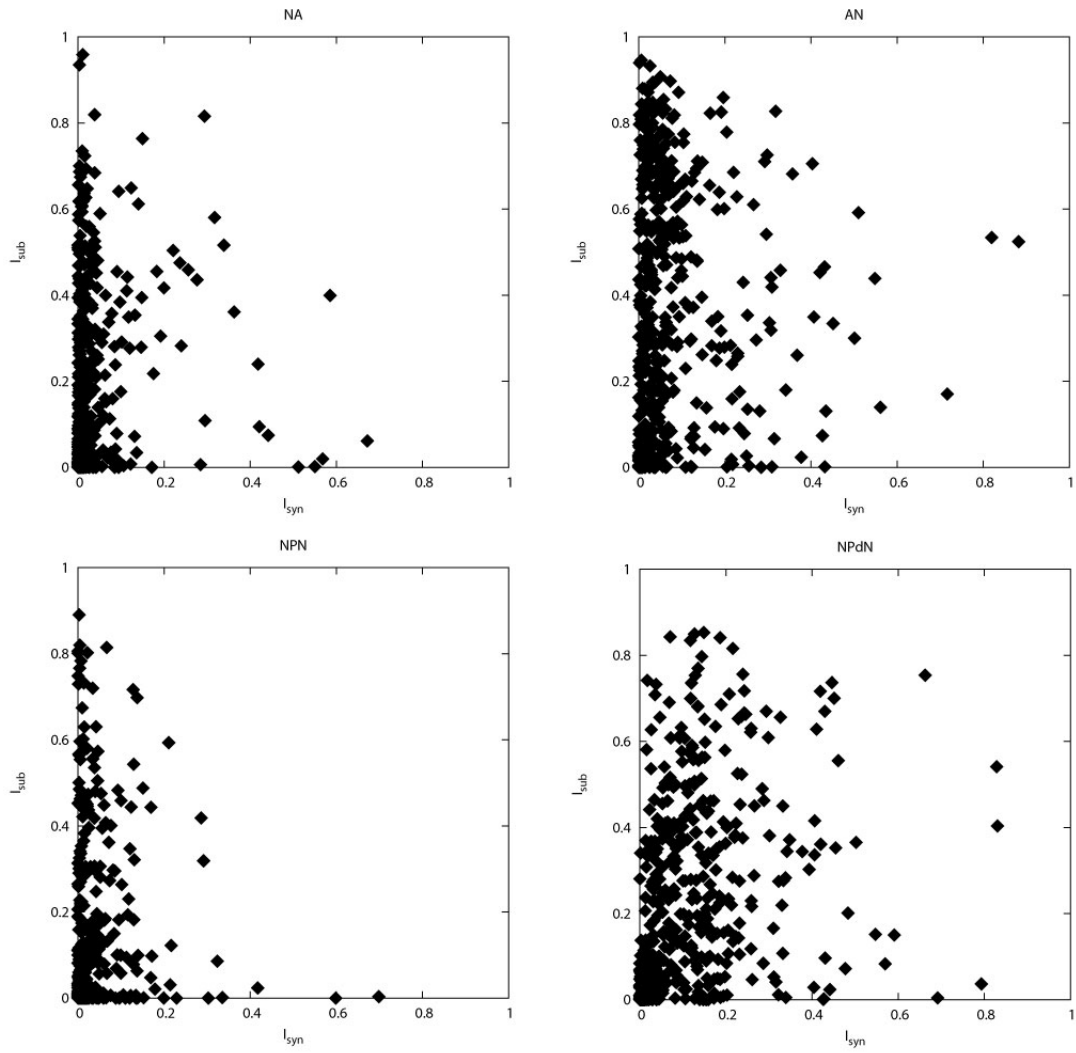


Figura 4.1: Distribuzione delle espressioni in input per i pattern NA, AN, NPN, NPdN a seconda dei valori empirici ottenuti per gli indici I_{syn} e I_{sub} .

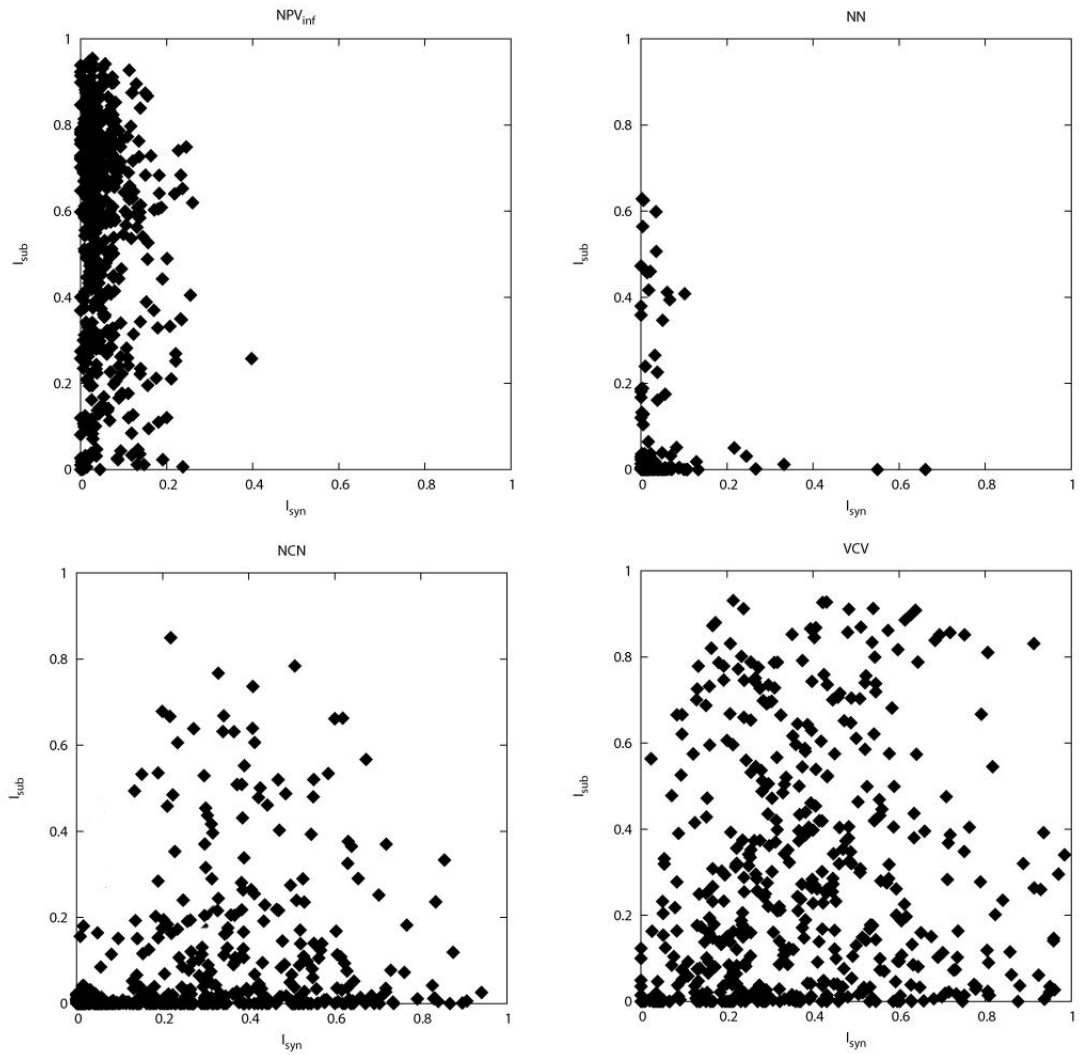


Figura 4.2: Distribuzione delle espressioni in input per i pattern NPV_{inf} , NN, NCN, VCV a seconda dei valori empirici ottenuti per gli indici I_{syn} e I_{sub} .

4.2.1 Analisi sul pattern NA

Il pattern NA risulta tipicamente associato alla formazione di sintagmi nominali, di cui il sostantivo ne è la testa sintattica. Come anticipato nel precedente capitolo, lo strumento computazionale è in grado di estrarre le espressioni più frequenti che corrispondano a tale pattern e, grazie ai tre livelli di annotazione del corpus, nel processo di selezione è in grado di individuare le espressioni che soddisfino le seguenti caratteristiche:

- **annotazione su lemma:** la combinazione deve essere formata da due elementi in sequenza, inclusi nella ricerca nella loro forma lemmatizzata;
- **annotazione PoS:** il primo elemento della combinazione deve essere contraddistinto dal tag specifico⁸ **S**, relativo alla categoria dei nomi comuni; il secondo deve essere etichettato con tag specifico **A**, relativo alla categoria degli aggettivi qualificativi;
- **annotazione sintattica:** l'aggettivo deve avere il sostantivo quale testa sintattica.

Una volta individuato l'insieme di tutte le espressioni che corrispondono ai criteri di cui sopra, tramite un ordinamento per occorrenze si trattengono le sole 500 espressioni più frequenti. Nello specifico, l'insieme di espressioni così selezionato spazia da valori di frequenza di 131.642 (*collegamento esterno*) a 735 (*mondo occidentale*); l'elenco completo delle espressioni ottenute è disponibile in Appendice B.

Si procede, quindi, all'esecuzione della procedura di test dello strumento, che genera per ogni espressione una serie di valori per ciascuno degli indici esposti nel precedente capitolo.

In riferimento all'indice di modificabilità sintagmatica, in questo caso vengono considerati sia il test di interrompibilità che quello di ordine inverso. In relazione al primo, l'interrompibilità è valutata in un numero libero di parole intervenienti entro il confine della frase, fintanto che il legame sintattico tra nome e aggettivo risulta preservato.

Per quanto specificato sull'indice di variabilità flessiva nel precedente paragrafo, a valle della generazione dei risultati è possibile considerare un piano individuato da due soli assi, rappresentanti la modificabilità sintagmatica e quella paradigmatica. Dalla disposizione delle espressioni in questo piano, è possibile individuare dei poli in relazione a valori massimi o minimi degli indici I_{syn} e I_{sub} , studiando le caratteristiche delle espressioni che vi si collocano.

La Tabella 4.4 mostra un sottoinsieme delle espressioni che esibiscono valori minimi per entrambi gli indici ($I_{syn}, I_{sub} < 0,01$) rappresentando, quindi, le

⁸Si veda l'appendice A per un elenco dei tag generali e specifici disponibili per il corpus.

Espr.	Freq.	I_{syn}	I_{sub}	Espr.	Freq.	I_{syn}	I_{sub}
trasporto pubblico	2.158	0,001579	0,009637	penisola iberico	809	0,003397	0,001235
campionato europeo	2.407	0	0,009057	sistema nervoso	1.649	0,007268	0,001211
emittente televisivo	969	0,000427	0,008188	atletica leggero	1.760	0,000243	0,001135
emisfero celeste	2.397	0,000763	0,007864	sistema operativo	6.416	0,005766	0,001090
energia cinetico	923	0,005816	0,007527	pietra miliare	953	0,005598	0,001048
elezione europeo	846	0,000955	0,007042	cemento armato	956	0	0,001045
cellula staminale	1.474	0	0,006069	testata giornalistico	1.001	0,000301	0,000998
equazione differenziale	825	0	0,006024	sci alpino	1.097	0	0,000911
oro olimpico	879	0,000348	0,005656	telefono cellulare	1.206	0	0,000829
campionato mondiale	3.760	0,005025	0,005554	campagna elettorale	2.691	0,004179	0,000743
stella gigante	1.173	0	0,005089	sito archeologico	2.816	0,001351	0,000710
torre campanario	941	0	0,004233	parlamento europeo	1.488	0,000005	0,000672
direttore artistico	1.218	0,007795	0,004088	galleria fotografico	3.082	0,000267	0,000649
carriera solista	1.054	0,004998	0,003781	etichetta discografico	1.780	0,000027	0,000561
politica estero	2.414	0,002632	0,003714	velocità radiale	2.220	0	0,000450
carro armato	3.224	0,001270	0,003708	colonna sonora	9.197	0,000503	0,000326
sistema solare	5.982	0,002513	0,002668	casa editore	5.148	0	0,000194
buco nero	1.583	0,000415	0,002520	navata centrale	1.243	0,008596	0
pallavolo femminile	834	0	0,002392	regione ecclesiastico	827	0,007589	0
magnitudine assoluto	2.139	0,002400	0,002332	gas naturale	863	0,003179	0
intelligenza artificiale	942	0,002665	0,002119	testamento biologico	1.126	0,002877	0
serie televisivo	9.057	0	0,002093	casa discografico	4.041	0,002810	0
tempo supplementare	964	0,003182	0,002070	occhio nudo	2.630	0,001344	0
scavo archeologico	980	0,000156	0,002037	classe spettrale	3.100	0,000119	0
stazione spaziale	991	0,007441	0,002014	cielo serale	2.603	0,000095	0
guerra freddo	1.643	0,000847	0,001823	circolo polare	1.583	0,000022	0
geografia fisico	1.652	0,000030	0,001813	amministratore delegato	1.082	0,000003	0
sala cinematografico	1.118	0,001526	0,001786	equatore celeste	2.050	0	0
strumento musicale	1.723	0,003525	0,001738	anidride carbonico	1.382	0	0
catena montuoso	1.571	0,006340	0,001271	miniserie televisivo	485	0	0

Tabella 4.4: Espressioni relative al pattern NA che presentano valori di variazione empirica inferiori all'1% per entrambi gli indici I_{syn} e I_{sub} . In rosso le unità di significato, in blu le espressioni riconducibili a terminologie di linguaggi tecnico-specialistici, in viola le espressioni il cui statuto è intermedio tra le due categorie o incerto.

entità che empiricamente si mostrano non modificabili né sintagmaticamente né paradigmaticamente⁹.

Analizzando le espressioni presenti sembra possibile individuare due principali categorie: espressioni che esibiscono una forte tendenza all'unitarietà semantica, come nei casi di *cemento armato*, *casa discografica*, *pietra miliare*, *strumento musicale*, ecc. (in rosso in Tabella 4.4) ed espressioni terminologiche proprie di linguaggi tecnico-specialistici (in blu), come nei casi di *anidride carbonica* (chimica), *equatore celeste*, *velocità radiale*, *cielo serale*, *stella gigante*, ... (astronomia), *sci alpino*, *pallavolo femminile*, ... (sport), *equazione differenziale* (matematica), *energia cinetica* (fisica), *cellule staminali* (medicina), *geografia fisica*, ... (geografia), *navata centrale* (arte), *amministratore delegato* (economia). Le espressioni in viola si collocano in gran parte a metà tra le due precedenti categorie¹⁰.

A queste espressioni si aggiunge un caso di combinazione cristallizzata in quanto parte di una locuzione avverbiale più ampia ([*ad*] *occhio nudo*).

⁹Da questo insieme sono state escluse le espressioni *quartier generale* e *valor militare* per la scorretta lemmatizzazione del primo componente in forma tronca.

¹⁰Nei casi di *sistema solare* o *buco nero*, ad esempio, risulta incerta la piena attribuzione delle espressioni al lessico specialistico dell'astronomia, in quanto esse sembrano largamente integrate nel linguaggio comune. Alcune espressioni come *parlamento europeo* o *carriera solista*, pur se integrate nel linguaggio comune, appaiono nel corpus in accezioni specialistiche del linguaggio giornalistico-giuridico e musicale rispettivamente.

Espr.	Freq.	I_{syn}	I_{sub}
conflitto mondiale	1.874	0,003191	0,935150
anno seguente	7.296	0,009638	0,735049
classe politico	847	0,003529	0,700281
operazione militare	1.120	0,003559	0,686645
crescita economico	853	0,005828	0,674800
età moderno	821	0,001217	0,656522
livello mondiale	1.180	0,002653	0,619014
classe sociale	1.134	0,002639	0,617393
lato sinistro	1.113	0,008905	0,607829
situazione economico	816	0,008505	0,605928
fonte storico	775	0,008951	0,594007
campagna militare	1.301	0,001535	0,587687
genere umano	910	0,002193	0,574643
braccio destro	1.298	0,004601	0,538650
centro cittadino	1.421	0	0,516168
azienda agricolo	777	0,007663	0,515898
scena musicale	768	0	0,510035
centro urbano	1.580	0,002525	0,507105
mano sinistro	943	0,001059	0,505402
situazione politico	901	0,003319	0,496213

Espr.	Freq.	I_{sub}	I_{syn}
adattamento cinematografico	785	0,001272	0,469165
livello internazionale	2.678	0,003349	0,466653
mano destro	1.244	0,000803	0,444240
diritto umano	3.081	0,002913	0,437834
corpo umano	1.312	0,009811	0,430981
fama internazionale	897	0,007743	0,428769
area metropolitano	1.473	0,001356	0,427757
cultura popolare	1.521	0,009766	0,412327
attività agonistico	960	0,009288	0,411577
mondo esterno	735	0,004065	0,403140
metro quadrato	1.213	0	0,394277
partito politico	3.142	0,002540	0,391831
risultato finale	1.216	0,007347	0,389533
parte finale	940	0,004237	0,384854
stile musicale	882	0,005637	0,352081
regime fascista	1.224	0,006494	0,348017
riva destro	744	0	0,344904
zona industriale	809	0,009792	0,341004
nucleo familiare	848	0	0,333296

Tabella 4.5: Espressioni relative al pattern NA che variazione sintagmatica inferiore all'1% ma variazione paradigmatica maggiore del 33%. In rosso le combinazioni i cui componenti mostrano tendenza alla riconoscibilità fraseologica, in blu le espressioni preferenziali meno coese.

In Tabella 4.5 sono invece mostrate le espressioni che hanno bassi valori di modificabilità sintagmatica (come sopra, inferiori all'1%), ma maggiore modificabilità paradigmatica ($I_{sub} > 0,33$).

Anche qui sembra possibile identificare due categorie di entità. Da un lato compaiono espressioni, in generale composizionali, caratterizzate da una certa riconoscibilità fraseologica, come *conflitto mondiale*, *crisi economica*, *genere umano*, ecc. In questo caso la non modificabilità sintagmatica si configura come la base di una loro coesione interna, nonostante esse siano espressioni a tutti gli effetti sostituibili (nel corpus si riscontra, infatti: *guerra* per *conflitto*; *crac*, *dissesto*, *recessione* per *crisi*; *specie* per *genere*). La sostituibilità, in particolare, sembra impedire a tali entità di procedere verso una cristallizzazione che le assimili a unità di significato.

Dall'altro lato sono presenti espressioni del tipo *giorno seguente*, *anno seguente*, *lato sinistro*, *mano sinistra*, ecc. che difficilmente sono inquadrabili come unità coese allo stesso modo delle precedenti, ma che risultano spesso inserite nei dizionari combinatori in quanto può essere loro riconosciuto un certo grado di familiarità tra i componenti¹¹.

In questo caso, a parità di comportamento empirico, il piano semantico sembra essere l'unico livello di analisi che permetta un discernimento tra tale gruppo e il precedente. In primo luogo è possibile osservare come nelle espressioni che abbiamo definito fraseologicamente riconoscibili l'aggettivo introduce nel sintagma una seconda entità tematica (dopo la testa) da cui esso deriva, assumendo un ruolo restrittivo¹². Per le espressioni meno coese del secondo gruppo (*mano sinistra*, *mondo esterno*), l'aggettivo non introduce alcuna nuova entità bensì assume un ruolo

¹¹Cfr., ad es., *mano destra/sinistra* in Lo Cascio (2011) e Tiberii (2012).

¹²Cfr. *crisi economica* = *crisi che riguarda l'economia*; *conflitto mondiale* = *conflitto che interessa il mondo*; *corpo umano* = *corpo dell'uomo*.

esclusivamente qualificativo e strettamente legato ad una proprietà della testa¹³. In secondo luogo, per le espressioni appartenenti al secondo gruppo l'aggettivo risulta tendenzialmente abbinato a tutta una serie di entità appartenenti ad una stessa classe: si pensi al caso di *giorno/settimana/mese/anno seguente* o di *lato/parte sinistra*. In questo senso si potrebbero identificare tali combinazioni come preferenze di selezione tra classi lessicali più che tra componenti, espandendo quindi il legame fraseologico su un piano più ampio¹⁴.

Anche tra le espressioni di Tabella 4.5, infine, compaiono frammenti di locuzioni avverbiali (*[a/di] livello mondiale*, *[a/di] livello internazionale*, *[di] fama internazionale*) che quindi non sono assimilabili alla categorizzazione operabile per i sintagmi nominali. Inoltre, due ulteriori espressioni meritano attenzione, a causa di peculiarità specifiche che le escludono dagli insiemi omogenei di cui si è detto sopra. *Braccio destro* ha un comportamento semantico particolare, in quanto tale espressione somma lettura letterale ed idiomatica al tempo stesso. Nel corpus, infatti, è riscontrabile un cospicuo numero di occorrenze in cui l'espressione viene utilizzata nel senso letterale di "parte del corpo" e quindi sostituibile, ad esempio da *arto destro* o anche *lato destro*, il che genera un valore medio-alto per I_{sub} ¹⁵. *Metro quadrato*, invece, chiaramente identificabile come un'unità di significato o, se si vuole, un'espressione terminologica, soffre della competizione con l'espressione sinonima sostituibile *metro quadro*. In questo caso, tuttavia, gli aggettivi *quadrato* e *quadro*, più che due lemmi indipendenti, possono essere visti come varianti dello stesso *item* lessicale e pertanto si può considerare, di fatto, l'intera espressione come non sostituibile.

La Tabella 4.6, infine, mostra tutte le espressioni con valori per l'indice di modificabilità sintagmatica maggiori del 20%. Come si vede dai dati, i casi di espressioni con valori medio-alti di variabilità sintagmatica ma bassa o nulla variabilità paradigmatica sono due e rientranti nel secondo dei casi visti nella tabella precedente, vale a dire delle espressioni preferenziali meno coese.

¹³Sarebbe possibile, in questo caso, richiamare l'idea di sintagma transitivo ed intransitivo di memoria modista, secondo cui non solo i verbi ma anche i sintagmi nominali possono veder transitare o meno il concetto da un elemento all'altro della costruzione (cfr. Graffi 2010, p.49). Nel caso delle espressioni del primo gruppo la transitività sarebbe possibile grazie alla seconda entità; in quelle del secondo gruppo, il concetto espresso dal sintagma ricadrebbe sulla testa e per questo verrebbe classificato come intransitivo.

¹⁴È interessante notare, in ogni caso, come gli aggettivi di tali combinazioni risultino tutti modificatori spaziali o temporali, ad indicazione di una possibile preferenza cognitiva che privilegi le delimitazioni di spazio e tempo nello stabilire abbinamenti preferenziali tra classi.

¹⁵Si ricordi, infatti, che lo strumento non opera una disambiguazione dei significati di una stessa entità lessicale. Se la semantica distribuzionale può aiutare a scartare le occorrenze di un'espressione sostituita che non è più sinonima dell'originale, non riesce altrettanto efficace a risolvere i casi in cui il vettore distribuzionale dell'espressione originale porta informazioni lessicali relative a due letture della stessa espressione. In una situazione di questo tipo, il massimo grado di sinonimia per un'espressione trasformata corrisponderebbe alla situazione in cui questa conservasse la doppia lettura a livello del significato. Poiché, in generale, quest'ultima eventualità non è verificata, per le espressioni trasformate che siano sinonime solo di uno dei significati originali ci si aspettano percentuali di sinonimia su livelli intermedi.

Espr.	Freq.	I_{syn}	I_{sub}
tempo stesso	2.772	0,671720	0,061403
famiglia nobile	1.192	0,585103	0,399281
settimana scorso	807	0,567987	0,019726
navata unico	747	0,511765	0,000959
anno scorso	2.956	0,441527	0,074559
tempo recente	1.135	0,417949	0,240005
azienda italiano	906	0,362421	0,361210
ruolo importante	2.352	0,338583	0,516237
personaggio famoso	750	0,316940	0,580529
caratteristica principale	1.312	0,295003	0,108705
periodo successivo	737	0,293384	0,816048
strada principale	793	0,276460	0,435973
società italiano	1.038	0,255914	0,459045
nome attuale	1.025	0,239614	0,282216
modo diverso	1.419	0,236686	0,475291
caratteristica peculiare	744	0,220942	0,504149

Tabella 4.6: Espressioni relative al pattern NA con variazione sintagmatica superiore al 20%. Il blu e il verde indicano rispettivamente le espressioni con valori bassi o medio-alti di modificabilità paradigmatica. In nero sono indicate le espressioni frammenti di locuzioni più ampie.

Per valori medi e alti, invece, di entrambi gli indici, si riscontrano espressioni che non sembrano esibire legami particolari tra i componenti (si veda *famiglia nobile*, *ruolo importante*, *azienda italiana*, ecc.) e che potremmo vedere come esempi di espressioni libere.

Anche in questo caso sono presenti frammenti di due espressioni avverbiali (*[al] tempo stesso* e *[in] tempi recenti*) e un'espressione aggettivale (*[a] navata unica*, in competizione nella trasformazione sintagmatica con il sintagma nominale *unica navata*).

Alla luce della distribuzione delle espressioni rispetto ai due indici di variazione, è quindi possibile proporre per le espressioni nominali relative al pattern NA una distinzione categoriale che vede l'intervento di tre etichette: consideriamo *polirematiche* le espressioni che presentano blocchi variazionali sia sintagmatici che paradigmatici (indipendentemente dalla composizionalità) e che comprendono sia entità unitarie nel significato, che espressioni terminologiche; chiameremo *collocazioni* quelle espressioni che presentano il solo blocco della variazione sintagmatica, formate da due elementi il cui contributo semantico è riconoscibile e generalmente legato alla presenza di due entità di cui una nel ruolo di modificatore della testa; definiamo infine *combinazioni preferenziali* le espressioni che presentano solo uno dei due blocchi rispetto alle variazioni sintagmatiche o paradigmatiche (preferenzialmente quello sintagmatico), il cui accostamento preferenziale è istituito tra una classe omogenea di entità nel ruolo di testa sintattica e un modificatore che non introduce nel sintagma una seconda entità.

La Tabella 4.7 mostra uno schema riassuntivo della categorizzazione proposta.

Modif. sintagmatica	Modif. paradigmatica	Categoria	Esempio
+	+	Espressione libera	<i>azienda italiana</i>
+	-	Combinazione preferenziale (?)	<i>anno scorso</i>
-	+	Collocazione	<i>crescita economica</i>
-	+	Combinazione preferenziale	<i>giorno seguente</i>
-	-	Polirematica	<i>sistema operativo</i>

Tabella 4.7: Categorizzazione delle espressioni nominali relative al pattern NA in relazione al comportamento empirico rispetto alle variazioni sintagmatiche e paradigmatiche. Il punto interrogativo identifica la categorizzazione incerta a causa della scarsità di espressioni nello specifico insieme.

È interessante notare, infine, come soffermandosi su un intervallo di valori molto stretto per l'indice I_{syn} , è possibile riconoscere che all'aumentare dei valori di I_{sub} le espressioni acquisiscano via via una coesione minore, coprendo il continuum da polirematica a collocazione e mostrando che tale categorizzazione risulta coerente non solo per i poli estremi considerati nello spazio $I_{syn}I_{sub}$, bensì anche per le espressioni che popolano la regione intermedia.

Espr.	Freq.	I_{syn}	I_{sub}
regime fascista	1.224	0,006494	0,348017
personaggio immaginario	1.027	0,006770	0,179259
edificio religioso	1.680	0,006505	0,176843
opera teatrale	2.264	0,006146	0,165484
comunità ebraico	1.165	0,006820	0,154048
tempo indeterminato	741	0,006702	0,131969
vita privato	3.536	0,006742	0,124658
centro commerciale	2.211	0,006292	0,022173
processo produttivo	880	0,006772	0,018478
mezzo pubblico	774	0,006418	0,011879
cellula staminale	1.474	0,006069	0
equazione differenziale	825	0,006024	0

Tabella 4.8: Esempio di variazione delle espressioni nel continuum categoriale all'aumento dell'indice I_{sub} , mantenendo I_{syn} fisso in una fascia ristretta di valori.

4.2.2 Analisi sul pattern AN

Analogamente a quanto visto per il pattern NA, anche per il pattern AN le 500 espressioni più frequenti vengono selezionate in base ai seguenti tre criteri:

- **annotazione su lemma:** la combinazione deve essere formata da due elementi in sequenza, inclusi nella ricerca nella loro forma lemmatizzata;
- **annotazione PoS:** il primo elemento della combinazione deve essere contraddistinto dal tag specifico **A**, relativo alla categoria degli aggettivi qualificativi;

il secondo deve essere etichettato con tag specifico **S**, relativo alla categoria dei sostantivi;

- **annotazione sintattica:** l'aggettivo deve avere il sostantivo quale testa sintattica.

Le espressioni così raccolte vanno da frequenza di occorrenza 31.110 (*maggior parte*, caso di errata lemmatizzazione del primo componente) a 489 (*pubblica sicurezza*) e sono disponibili in Appendice C. Anche in questo caso la categoria grammaticale in uscita, quella cioè dell'intero sintagma, è di norma nominale.

Il pattern AN ha la peculiarità di corrispondere ad una combinazione che, relativamente ai sintagmi nominali, risulta marcata rispetto alla sequenza non marcata NA. In generale, infatti, la sintassi italiana vede la collocazione dei modificatori a destra della testa: mentre per gli aggettivi che prevedono la reggenza di un complemento tale configurazione è obbligatoria (es. *uomo ligio al dovere*), una eventuale dislocazione a sinistra dell'aggettivo si deve, invece, a particolari ragioni stilistiche o semantiche. È noto¹⁶, infatti, che tanto più la qualità espressa dall'aggettivo è oggettiva, e quindi relativa ad uno stato fisico, tanto più la sua posizione risulta fissa e post-nominale (ad es. *segnaletica stradale* con 90 occorrenze in PAISÀ, ma **stradale segnaletica*, non attestata). Quando invece la qualificazione espressa dall'aggettivo è soggettiva e legata alle considerazioni del parlante, essa permette la dislocazione a sinistra del modificatore, investendo spesso quest'ultimo di nuove sfumature di significato meno legate all'accezione puramente "fisica" (cfr. *amico vecchio* vs. *vecchio amico*).

Tale peculiarità gioca un ruolo essenziale nell'analisi dei risultati di tale pattern, condizionando il test di modificabilità sintattica, come si vedrà a breve, che anche in questo caso si esplica attraverso la ricerca di espressioni interrotte o i cui componenti siano in ordine inverso.

La Tabella 4.9 mostra l'insieme di espressioni che esibiscono un blocco di modificabilità sintagmatica e paradigmatica ($I_{syn} < 0,01$, $I_{sub} < 0,1$ ¹⁷) una volta che le espressioni ripetute, quelle per cui i tag presenti nel corpus fossero errati¹⁸ o quelle in cui la lemmatizzazione risultava errata siano state escluse. Come si vede, le espressioni sono principalmente di due tipologie: espressioni fortemente cristallizzate entrate nell'uso comune ormai come unità di significato (*nuovo millennio*, *alta velocità*, *grande schermo*) e frammenti di locuzioni avverbiali che risultano cristallizzate solo in relazione all'espressione più ampia che li ingloba (*[ad un] certo punto*,

¹⁶Per quanto non siano molti gli studi che hanno affrontato tale interessante fenomeno dell'italiano, per una visione più completa si rimanda ai contributi di Alisova (1967), D'Addio (1974), Serianni (1989, pp. 199-205), Vincent (1986).

¹⁷In questo caso si è considerato un intervallo più ampio per l'indice I_{sub} per includere un maggior numero di espressioni nell'insieme.

¹⁸Nel corpus, ad esempio, le espressioni *comune italiano*, *fine stagione* o *lato opposto* sono etichettate come combinazioni AN.

[a] pari merito, [nelle] immediate vicinanze, [un] bel po' [di], [nella/la] stragrande maggioranza [di], ecc.). Si noti la presenza di *grosso modo* come unica costruzione avverbiale effettivamente assimilabile al pattern AN non in quanto frammento di locuzione.

Espr.	Freq.	I_{syn}	I_{sub}
alto velocità	2.806	0,009880	0,091472
grosso modo	591	0,001689	0,046549
massimo campionato	733	0,001362	0,033323
pari opportunità	552	0	0,023116
lieto fine	585	0,005102	0,012066
massimo serie	3.003	0	0,002850
buono domenica	684	0	0,001480
libero arbitrio	505	0	0,000341
pronto soccorso	1.034	0,000966	0,000030
social network	1.124	0,007944	0
personal computer	913	0,001094	0
nuovo millennio	692	0	0

Espr.	Freq.	I_{syn}	I_{sub}
certo punto	3.183	0,001568	0,077479
pari merito	607	0	0,059655
immediato vicinanza	733	0,001362	0,059388
massimo splendore	532	0	0,058243
lungo termine	1.576	0,005051	0,006693
stesso anno	23.160	0,004256	0,000768
tenero età	869	0,002296	0,000607
stesso sesso	594	0,005025	0,000515
lungo raggio	1.101	0,009001	0,000255
pari passo	686	0,008671	0,038345
bello po'	592	0,008375	0
stragrande maggioranza	1.230	0,002433	0

Tabella 4.9: Espressioni relative al pattern AN che presentano valori di variazione empirica inferiori all'1% per l'indice I_{syn} e inferiori al 10% per I_{sub} . A sinistra sono mostrati i sintagmi nominali; a destra frammenti di locuzioni.

In Tabella 4.10 sono invece mostrate le espressioni che esibiscono bassi valori di modificabilità sintagmatica ma alta sostituibilità. La maggior parte di queste sono,

Espr.	Freq.	I_{syn}	I_{sub}
grande popolarità	516	0,005780	0,945849
alto carica	566	0,008757	0,880606
grande influenza	615	0,006462	0,843848
massimo divisione	669	0,001493	0,818857
nuovo costituzione	550	0,005425	0,808753
intero area	581	0,006838	0,590250
grande potenza	1.028	0,004840	0,588217
ex marito	628	0,001590	0,576427
intero territorio	821	0,002430	0,576088
grande maestro	912	0,007617	0,562590
mezzo secolo	549	0,001818	0,508773
nuovo imperatore	501	0,009881	0,506071
ultimo volta	2.502	0,009501	0,383282

Espr.	Freq.	I_{syn}	I_{sub}
grande interesse	1.293	0,001544	0,796725
grande importanza	2.192	0,003636	0,759837
grande attenzione	604	0,009836	0,730186
certo importanza	620	0,003215	0,726053
ultimo momento	1.008	0,008850	0,689801
notevole interesse	770	0,006452	0,669859
vario genere	1.056	0,009381	0,656953
grande valore	1.134	0,008741	0,646957
vario tipo	1.282	0,006971	0,625059
grande varietà	1.157	0,004303	0,564404
stesso periodo	4.671	0,002349	0,549183
poco distanza	1.056	0	0,507810
grande difficoltà	503	0,005929	0,466194
grande capacità	788	0,002532	0,438152
buon livello	492	0	0,436570
grande quantità	2.719	0,004758	0,434845
scarso successo	922	0,004320	0,424828
largo scala	1.110	0,000900	0,387148
alto definizione	898	0,002222	0,377096
nuovo governo	1.570	0,004439	0,374653
lungo durata	684	0,004367	0,371636
nuovo generazione	2.720	0,009829	0,368140

Tabella 4.10: Espressioni relative al pattern AN che presentano valori di variazione empirica inferiori all'1% per l'indice I_{syn} e superiori al 33 % per I_{sub} . A sinistra i sintagmi nominali, a destra i frammenti di locuzioni più ampie.

ancora una volta, frammenti di locuzioni più ampie, come mostrato nella seconda delle Tabelle 4.10. Per quanto riguarda, invece, le espressioni che costituiscono sintagmi nominali a tutti gli effetti (tabella di sinistra), si verifica che in generale per tutte le espressioni il numero di sostituzioni è molto maggiore per il secondo componente (il nome) piuttosto che per il modificatore (cfr. Tabella 4.11). Tale evidenza è in realtà conseguenza di due fenomeni distinti.

Una parte delle espressioni presenta un aggettivo che, in italiano, ha l'obbligo o una forte tendenza a presentarsi in posizione prenominali (*ultimo, mezzo, ex*). In

Espr.	Sin2/Sin1
grande popolarità	151,9
alto carica	16,2
grande influenza	22,1
massimo divisione	223
nuovo costituzione	39,8
intero area	37,3
grande potenza	88
ex marito	854
intero territorio	62,1
grande maestro	8,4
mezzo secolo	1311,5
nuovo imperatore	4,8
ultimo volta	1800,5

Tabella 4.11: Rapporto fra il numero di occorrenze ottenute tramite sostituzione del secondo componente (N) rispetto a quelle ottenute sostituendo il primo (A) per i sintagmi nominali del pattern AN che esibiscono sostituibilità empirica ma non modificabilità sintagmatica.

queste situazioni, eccetto eventuali pochi casi di interruzione, la modificabilità sintagmatica è inibita rispetto all'inversione dei componenti. La sostituzione, inoltre, può avvenire solo per il nome, in quanto, nonostante sia potenzialmente possibile anche per l'aggettivo, questo possiede di norma sinonimi che se sostituiti seguirebbero il nome nei sintagmi non marcati, invece di precederlo. Si pensi, a riguardo, al caso di *ultima volta*, in cui i sinonimi dell'aggettivo quali *decisivo*, *definitivo* appaiono, di norma, a destra di *volta*¹⁹. Le espressioni che rientrano in questo insieme non possono, quindi, essere viste come entità di interesse fraseologico, in quanto i valori empirici di cristallizzazione sintagmatica sono in parte dovuti a restrizioni sintattiche intrinseche della sintassi italiana, in parte ai limiti dello strumento.

Le rimanenti espressioni a sinistra in Tabella 4.10 mostrano anch'esse una prevalenza di sostituzione per il sostantivo, e sono accomunate dalla presenza dei soli aggettivi *grande*, *nuovo*, *intero*, *alto* in posizione prenominale. Tali aggettivi mostrano una spiccata tendenza ad assumere accezioni diverse a seconda della loro posizione precedente o seguente il nome²⁰. In questi casi si vede come il nome, o la classe di entità cui esso appartiene, è abituato a selezionare semanticamente solo l'accezione aggettivale relativa alla posizione prenominale, impedendo gli accostamenti con lo stesso aggettivo dislocato a destra (cfr. *nuovo imperatore*, 501 occorrenze, con **imperatore nuovo*, non attestato nel corpus). Anche in questo caso, quindi, la non modificabilità sintagmatica non sembra essere legata a preferenzialità di carattere

¹⁹In PAISÀ *volta decisiva* ottiene 5 occorrenze, *volta definitiva* 39.

²⁰*Grande* e *alto* assumono il significato di *importante* in posizione prenominale invece dei rispettivi significati legati alle dimensioni fisiche assunti quando in posizione postnominale; *nuovo* prenominale è legato all'idea di recente cambiamento o entrata in vigore, piuttosto che al significato postnominale di "recente creazione"; *intero*, quando preposto al nome, identifica il significato di "totale", più che del significato postnominale tipico di "integro" o "intatto".

lessicale, bensì ad una proprietà di selezione semantica tipica della lingua.

In Tabella 4.12 sono mostrate, invece, le espressioni che presentano valori medi o alti di variazione sintagmatica ($I_{syn} > 0,2$) e il blocco delle sostituzioni. Ancora

Espr.	Freq.	I_{syn}	I_{sub}
medio dimensione	526	0,426390	0,073801
omonimo film	640	0,377432	0,023559
alto voce	523	0,314548	0,066796
omonimo romanzo	809	0,309137	0,001199
pubblico amministrazione	1.416	0,284125	0,000175
pubblico sicurezza	489	0,255708	0,003766
ultimo anno	14.021	0,233280	0,091446
prossimo mese	615	0,217557	0,006706
prossimo settimana	863	0,216878	0,008286
pubblico istruzione	705	0,214922	0,018967
ultimo mese	1.450	0,204608	0,001214

Tabella 4.12: Espressioni relative al pattern AN che presentano valori per l'indice I_{syn} maggiori di 0,2 e valori per l'indice I_{sub} minori di 0,1.

una volta compaiono frammenti di locuzioni avverbiali ($[di]$ *medie dimensioni*, $[ad]$ *alta voce*), che in questo caso ammettono l'inversione dei costituenti ($[di]$ *dimensioni medie*, $[a]$ *voce alta*).

Anche le espressioni temporali *ultimo anno*, *prossimo mese*, *prossima settimana*, *ultimo mese*, a differenza di quanto accadeva per espressioni analoghe nel pattern NA, hanno una forte tendenza a comparire in locuzioni (*negli ultimi anni*, *nei prossimi mesi*, ecc.) venendo spesso interrotte da un quantificatore numerico (*negli ultimi NUM anni*, *nelle prossime NUM settimane*). Le tre espressioni in cui compare l'aggettivo *pubblico* costituiscono un caso particolare e devono la loro modificabilità sintagmatica alla competizione tra forme cristallizzate proprie del linguaggio giuridico-politico e forme sciolte incluse nella sintassi libera con aggettivo posposto²¹.

Per quanto riguarda, infine, le espressioni in cui compare l'aggettivo *omonimo*, esse sembrano le sole ad esibire una qualche forma di preferenzialità combinatoria in quanto, nonostante *omonimo* non abbia sinonimi esprimibili attraverso un solo lemma, esso viene associato a *romanzo* o *film* (sostantivi che ammettono sinonimi quali *racconto*, *libro* o *pellicola*) in proporzioni maggiori.

La Tabella 4.13, infine, mostra le espressioni con valori medio alti per entrambi gli indici ($I_{syn} > 0,3$; $I_{sub} > 0,33$).

²¹Si confrontino i seguenti esempi tratti dal corpus:

Fra il pubblico c'era anche Tullio De Mauro, ex ministro della *Pubblica Istruzione*.
Promosse l'*istruzione pubblica*, creando alcuni collegi e licei statali.

Espr.	Freq.	I_{syn}	I_{sub}
solo anno	595	0,882271	0,524336
importante ruolo	522	0,819440	0,534196
principale attività	492	0,548209	0,438762
maggiore dimensione	559	0,510079	0,592032
prossimo anno	1.410	0,451362	0,334141
migliore risultato	517	0,431243	0,465980
solo punto	567	0,420838	0,452703
principale centro	849	0,406709	0,349377
diverso tipo	783	0,403201	0,705477
elevato numero	688	0,357009	0,681607
stesso tempo	5.668	0,328675	0,457889
ultimo secolo	500	0,317872	0,827412
recente studio	760	0,309091	0,419053
fondamentale importanza	648	0,306952	0,441037
principale città	1.061	0,303349	0,336157

Tabella 4.13: Espressioni relative al pattern AN che presentano valori per l'indice I_{syn} maggiori di 0,3 e valori per l'indice I_{sub} maggiori di 0,33.

Come si vede, oltre alle espressioni che rappresentano frammenti di locuzioni più ampie (*[allo] stesso tempo*), non sembra riscontrabile alcuna espressione che mostri marcatamente un qualche legame fraseologico tra i componenti.

Alla luce di quanto esposto, e in virtù del valore marcato del sintagma AN in italiano, sembra possibile identificare chiaramente il solo insieme delle polirematiche in relazione alle espressioni che hanno un blocco variazionale sia sintagmatico che paradigmatico e che non appartengano a locuzioni che le inglobino. Per gli altri poli di variabilità, la grande disomogeneità di comportamenti e le interferenze con le proprietà sintattiche e semantiche del pattern non sembrano favorire l'identificazione univoca e affidabile di alcuna categoria fraseologica.

4.2.3 Analisi sul pattern NPN

Il pattern NPN si configura, ancora una volta, come una sequenza grammaticale che genera espressioni nominali, formate da una testa (il primo dei sostantivi dell'espressione) e un sintagma preposizionale composto dalla preposizione e il secondo sostantivo. I criteri di selezione delle espressioni sono i seguenti:

- **annotazione su lemma:** la combinazione deve essere formata da tre elementi in sequenza, inclusi nella ricerca nella loro forma lemmatizzata;
- **annotazione PoS:** il primo ed il terzo elemento della combinazione devono essere contraddistinti dal tag specifico **S**, relativo alla categoria dei sostantivi; il secondo deve essere etichettato con il tag specifico **E**, relativo alla categoria delle preposizioni semplici;

- **annotazione sintattica:** la preposizione deve avere il sostantivo quale testa sintattica, mentre il secondo sostantivo deve dipendere dalla preposizione.

Le espressioni così selezionate spaziano da un massimo di occorrenze pari a 24.259 (*punto di vista*) fino ad un minimo di 352 (*ora di lavoro*), e sono riportate in Appendice D. Per questo pattern lo strumento opera unicamente il test di interrompibilità nell'ambito della verifica della modificabilità sintagmatica, in quanto l'inversione di ordine dei due sostantivi comporterebbe la creazione di espressioni totalmente differenti in significato da quella originale oppure prive di senso (cfr. *casa di cura* → *cura di casa* o *polvere da sparo* → **sparo da polvere*).

Analogamente a quanto visto per le combinazioni NA, anche in questo caso il blocco di modificabilità sintagmatica e paradigmatica è caratteristico di espressioni che costituiscono unità di significato oppure sequenze fortemente terminologizzate, come è possibile vedere in Tabella 4.14, dove compaiono espressioni come *polvere da sparo*, *ferro di cavallo*, *canna da zucchero*, *colpo di stato* relativamente al primo gruppo e *pallone d'oro*, *funzione d'onda* in relazione al secondo.

In questo insieme compaiono, tuttavia, anche espressioni composte da coppie di quantificatori (*centinaia di migliaia*, *decine di migliaia*), costruzioni in cui è presente un quantificatore ed un'unità di tempo, misura o valuta (*paio di mesi*, *milione di tonnellate*, *centinaia di metri*, *milioni di dollari*, *miliardi di euro*), espressioni denominative del tipo *mese di* (*mese di febbraio/ottobre/novembre/ecc.*) e costruzioni terminologiche che coinvolgono tassi (*[NUM] valvole per cilindro*, *[NUM] euro a persona*). La presenza di tali espressioni tra le entità non modificabili è motivato da due principali fattori: da un lato tali combinazioni non sono in genere interrompibili per la loro natura di costruzioni di quantificazione e/o denominazione; in secondo luogo la sostituibilità è inibita in quanto i lemmi che li compongono sono principalmente terminologici o numerici e pertanto privi di sinonimi²².

Si segnalano, infine, nell'insieme, frammenti di locuzioni avverbiali o aggettivali, come *[in] fin di vita*, *[a] portata di mano*, *[di] volta in volta*, *[senza] soluzione di continuità*, *[al] giorno d'oggi*, *[senza] ombra di dubbio*, oltre al frammento del sintagma nominale più ampio *visibilità ad occhio [nudo]*. Esiste anche la locuzione avverbiale *mano a mano*, corrispondente esattamente al pattern NPN.

La Tabella 4.15 mostra invece le espressioni caratterizzate da bassi valori di interrompibilità ma alti valori di sostituibilità.

Come si vede, in questo insieme compaiono due tipologie di entità. In primo luogo, espressioni i cui componenti mantengono gran parte della propria indipendenza semantica, ma esibiscono una certa preferenzialità di combinazione lessicale (*libertà di espressione*, *via di fuga*, *porta d'ingresso*, *luogo di lavoro*, ecc.). In secon-

²²Tale risultato pone in risalto un limite importante dello strumento che, per alcune specifiche classi omogenee di entità (ad es. sostantivi numerici come *migliaia*, *milioni*, ecc.; unità di tempo, misura e valuta; nomi di mesi e stagioni) potrebbe prendere in considerazione, per il test di sostituibilità, non solo i sinonimi (peraltro inesistenti per le classi suddette), bensì tutti gli elementi appartenenti alla classe, come già operato in Lin (1999).

Espr.	Freq.	I_{syn}	I_{sub}
disco d' oro	1.005	0	0,009538
istituto di credito	612	0,006494	0,009524
campo di battaglia	2.137	0,001868	0,008264
uomo d' affare	1.099	0,003626	0,007852
segretario di stato	392	0	0,007842
casa di cura	354	0	0,007795
casa di produzione	1.474	0,005398	0,007691
cadenza di tiro	553	0,005396	0,007134
decina di anno	494	0,008032	0,006826
libertà di stampa	758	0,003942	0,006821
proposta di legge	926	0,003229	0,006805
miliardo di dollaro	2.506	0,005556	0,006794
ombra di dubbio	461	0,006466	0,006626
amico di famiglia	362	0,008219	0,006570
densità di popolazione	943	0,008412	0,006535
corpo di fabbrica	571	0,005226	0,006080
campo di concentramento	2.358	0	0,006050
modalità di gioco	906	0,007667	0,005602
arma da fuoco	1.734	0,004021	0,005493
guerra di successione	924	0,002160	0,005106
edificio di culto	1.008	0,004936	0,004883
gioco d' azzardo	752	0,003974	0,004826
gioco di ruolo	1.399	0,002851	0,004792
decina di metro	525	0,003795	0,004776
battuta di caccia	435	0	0,004437
messa in scena	899	0	0,004411
paio di giorno	409	0,002439	0,003605
messa a punto	702	0	0,003503
casa di riposo	479	0,004158	0,003180
crimine di guerra	649	0,001538	0,002994
rapporto di compressione	864	0,002309	0,002895
forza di gravità	555	0,001799	0,002649
marchio di fabbrica	455	0,006550	0,002549
punto di svolta	516	0,005780	0,002493
opera d' arte	3.367	0,003551	0,002315
anno di regno	597	0,008306	0,002249
impianto di risalita	414	0,009569	0,002228
parola d' ordine	608	0,001642	0,002030
metro di diametro	400	0,007444	0,001945
giorno d' oggi	1.555	0	0,001921
sistema d' arma	598	0,006645	0,001878
salto di qualità	739	0,001351	0,001869
mal di testa	414	0	0,001836
guerra d' indipendenza	902	0,001107	0,001722
centinaio di metro	1.091	0,000916	0,001398
testa di serie	389	0,002564	0,001352
presa di posizione	732	0	0,001319
girone di ritorno	686	0	0,001292
visibilità ad occhio	756	0	0,001263
figlio d' arte	730	0,001368	0,001252
scontro a fuoco	428	0	0,001169
colpo di testa	427	0	0,001112
testa di ponte	442	0,002257	0,000922
nome in codice	1.411	0,004234	0,000903
compagno di squadra	2.733	0,000366	0,000895
carriera di allenatore	362	0,005495	0,000837
guerra di indipendenza	1.205	0,005776	0,000782
presa d' aria	627	0,001592	0,000689
bocca da fuoco	374	0	0,000633
lunghezza d' onda	1.939	0	0,000614
miliardo di euro	2.132	0,003738	0,000521
periodo di visibilità	1.272	0,006250	0,000493
anno di piombo	490	0,006085	0,000492
fibra di carbonio	707	0,001412	0,000473
via di mezzo	495	0,006024	0,000461
mano a mano	549	0	0,000453
onda d' urto	465	0	0,000427
titolo di coda	585	0	0,000418
uscita di scena	437	0,002283	0,000415

Espr.	Freq.	I_{syn}	I_{sub}
motore di ricerca	1.183	0,000845	0,000403
sedia a rotella	683	0	0,000398
borsa di studio	1.605	0	0,000370
metro di profondità	697	0,009943	0,000367
soluzione di continuità	535	0	0,000367
pò di tempo	1.242	0	0,000320
guerra di secessione	673	0,001484	0,000320
diritto d' autore	1.717	0,004638	0,000299
milione di dollaro	6.120	0,007782	0,000261
agente di polizia	562	0,003546	0,000230
casa di moda	373	0,007979	0,000186
punto di vista	24.159	0,000124	0,000062
secondo d' arco	748	0	0,000036
euro a persona	537	0,009225	0
milione di tonnellata	702	0,008475	0
mese di luglio	938	0,008457	0
sci di fondo	492	0,008065	0
valvola per cilindro	556	0,007143	0
mese di agosto	841	0,005910	0
mese di aprile	675	0,005891	0
mese di dicembre	535	0,005576	0
obiezione di coscienza	385	0,005168	0
mese di giugno	792	0,005025	0
volta in volta	452	0,004405	0
disco di platino	1.209	0,004119	0
prigioniero di guerra	741	0,004032	0
mese di settembre	755	0,003958	0
atomo di carbonio	523	0,003810	0
minuto per minuto	547	0,003643	0
maestro di cappella	561	0,003552	0
condanna a morte	842	0,003550	0
decina di migliaio	1.206	0,003306	0
vicino di casa	354	0,002817	0
mese di maggio	853	0,002339	0
portata di mano	440	0,002268	0
gioco di parola	920	0,002169	0
posto a sedere	513	0,001946	0
numero di telefono	523	0,001908	0
mese di novembre	568	0,001757	0
direttore d' orchestra	1.235	0,001617	0
comandante in capo	768	0,001300	0
mese di ottobre	838	0,001192	0
olio su tela	919	0,001087	0
centinaio di migliaio	1.191	0,000839	0
datore di lavoro	1.711	0,000584	0
messa in onda	1.743	0,000573	0
colpo di stato	2.183	0,000458	0
pena di morte	1.505	0	0
calcio di rigore	954	0	0
cavallo di battaglia	599	0	0
fin di vita	586	0	0
fuoco d' artificio	547	0	0
colpo di grazia	537	0	0
specchio d' acqua	534	0	0
funzione d' onda	517	0	0
permesso di soggiorno	509	0	0
canna da zucchero	499	0	0
cacciatore di taglia	495	0	0
carrello d' atterraggio	491	0	0
ferro di cavallo	443	0	0
posto di blocco	428	0	0
mese di febbraio	426	0	0
polvere da sparo	414	0	0
ponte di volo	394	0	0
pallone d' oro	385	0	0
corpo a corpo	363	0	0
paio di mese	361	0	0
girone d' andata	356	0	0

Tabella 4.14: Espressioni relative al pattern NPN che presentano valori di variazione empirica inferiori all'1% per l'indice I_{syn} e I_{sub} .

Espr.	Freq.	I_{syn}	I_{sub}
tenore di vita	389	0,002564	0,890432
progetto di legge	459	0,004338	0,819910
anno d' età	488	0	0,806194
luogo di lavoro	565	0	0,800924
metro di altitudine	451	0,006608	0,783475
album d' esordio	534	0,003731	0,766233
anno di prigione	359	0	0,748530
lasso di tempo	908	0,001100	0,729616
metro d' altezza	411	0,009639	0,674383
terreno di gioco	442	0,004505	0,596527
porta d' ingresso	519	0,009542	0,573211
via d' uscita	489	0,002041	0,566513
metro di quota	387	0,002577	0,564449
punto di incontro	374	0,005319	0,554133
anno di età	1.735	0,002300	0,500793
via di fuga	629	0,003170	0,485244
linea di confine	401	0,007426	0,471046
mezzo di produzione	370	0,005376	0,462406
paese d' origine	525	0	0,452596
periodo d' oro	488	0,006110	0,445827
ciclo di vita	388	0,007673	0,353324
libertà di espressione	649	0,004601	0,341373
anno di vita	2.304	0,006040	0,336166

Tabella 4.15: Espressioni relative al pattern NPN che presentano valori di variazione empirica inferiori all'1% per l'indice I_{syn} e superiori al 33% per l'indice I_{sub} .

do luogo, è facile riconoscere anche qui costruzioni di quantificazione come *[NUM] anni di vita/di età/di prigionia, [NUM] metri di quota/d'altezza/d'altitudine*. Queste espressioni mostrano inibizione alla modificabilità sintagmatica a causa della fissità della costruzione quantificativa *[NUM N1 di N2]*, tuttavia il componente N2 appare liberamente sostituibile.

In Tabella 4.16 sono mostrate, invece, le espressioni che esibiscono valori empirici di modificazione sintagmatica medio-alti ($I_{syn} > 17\%$).

Espr.	Freq.	I_{syn}	I_{sub}
campionato di calcio	1.659	0,598597	0
causa di problema	725	0,417203	0,023550
campione in carica	635	0,335079	0,001204
famiglia di origine	694	0,323587	0,085734
olio di oliva	369	0,302457	0
serie di evento	464	0,291603	0,319030
attacco da parte	520	0,285714	0,418520
campionato di serie	885	0,229094	0
ragazzo di nome	395	0,216270	0,122286
meta di pellegrinaggio	371	0,213983	0,031019
numero di persona	828	0,210677	0,593458
olio d' oliva	503	0,199045	0
successo di critica	423	0,178641	0,020771

Tabella 4.16: Espressioni relative al pattern NPN che presentano valori di variazione sintagmatica maggiori del 17%.

Si vede che le espressioni, in tale insieme, che mostrano valori minimi di sostituzione (*olio di/d' oliva, campionato di calcio, campione in carica, famiglia di origine, successo di critica*) sono entità che esibiscono una preferenzialità di coselezione lessicale ed una riconoscibilità fraseologica.

Le espressioni *causa di problemi, campionato di serie, attacco da parte* e *ragazzo di nome* risultano frammenti di sintagmi più ampi intercettati dal pattern in questione e, specificamente, parti delle locuzioni *[a] causa di problemi [A], campionato di serie [N], attacco da parte di [N]* e *ragazzo di nome [N_{pr}]*.

Le espressioni che, invece, mostrano valori medio-alti per entrambi gli indici (*serie di eventi* e *numero di persone*) risultano espressioni quantificative che non sembrano esibire particolare attrazione tra i componenti.

Alla luce di quanto visto è possibile anche qui riassumere la categorizzazione delle espressioni in base al loro comportamento variazionale come segue: le espressioni che non esibiscono variabilità né per interruzione né per sostituzione possono essere in generale considerate *polirematiche* quando non siano parte di un costrutto denominativo o quantificativo. Riguardo a quest'ultima tipologia, è interessante notare che Masini (2007, p. 166) include le costruzioni *[N1 di]* e *[NUM N1 di]*, parti delle più ampie *[[N1 di] N2]* e *[[NUM N1 di] N2]*, nell'insieme delle espressioni multiparola modificatori di elementi nominali, sotto il nome di *delimitatori*. Nel nostro caso, date le espressioni presenti nell'insieme, è possibile generalizzare il pattern

dei delimitatori con modificatore numerico alla sequenza NUM N1 P N2, in modo da includere anche le costruzioni con preposizione *a* o *per* (*NUM euro a persona*, *NUM valvole per cilindro*). Vale, per queste costruzioni, il principio per cui è la sequenza grammaticale in sé (costruzione) a generare il valore delimitativo, più che la scelta lessicale, in quanto il sostantivo N2 appare libero e N1 solo parzialmente condizionato²³.

Le espressioni, invece, che esibiscono in maniera lineare solo uno dei blocchi variazionali, sia esso sintagmatico o paradigmatico, sembrano esibire una riconoscibilità fraseologica tale da poter essere definite *collocazioni*.

Infine, non ci sono abbastanza dati a disposizione per inferire, al pari del pattern NA, che le espressioni ad alta variazione sia sintagmatica che paradigmatica siano espressioni libere. Le uniche due espressioni con tali caratteristiche (*serie di eventi* e *numero di persone*) apparirebbero ancora come sequenze quantificative, benché esibiscano sia un alto numero di sostituzioni (per il primo e per il secondo sostantivo), sia interruzioni (tra primo e secondo componente e tra secondo e terzo). Non va esclusa la possibilità che la sequenza grammaticale NPN, in generale, risulti un pattern non incluso a pieno titolo nella sintassi libera dell'italiano per la mancanza di un determinante interno (a differenza del pattern NPdN, analizzato di seguito). Come ricorda Masini (2008, p.567): «L'assenza del determinante, e quindi del principale "contorno sintattico" del nome [...] rimanda a fenomeni quali l'incorporazione o la composizione, in cui gli elementi coinvolti perdono forza referenziale e libertà sintattica in quanto coinvolti in operazioni lessicali o pseudo-lessicali».

La seguente Tabella 4.17 schematizza i risultati per il pattern NPN.

Modif. sintagmatica	Modif. paradigmatica	Categoria	Esempio
+	+	Espr. libera/Costr. quantificativa (?)	<i>serie di eventi</i>
+	-	Collocazione	<i>olio d'oliva</i>
-	+	Collocazione	<i>via di fuga</i>
-	-	Polirematica	<i>fuoco d'artificio</i>
-	-	Costruzione quantificativa	<i>paio di mesi</i>
-	-	Costruzione denominativa	<i>mese di febbraio</i>

Tabella 4.17: Categorizzazione delle espressioni nominali relative al pattern NPN in relazione al comportamento empirico rispetto alle variazioni sintagmatiche e paradigmatiche.

4.2.4 Analisi sul pattern NPdN

Il pattern NPdN è una variazione del pattern analizzato nel precedente paragrafo, in cui la preposizione viene scelta articolata. Anche in questo caso la sequenza genera, di norma, espressioni nominali.

I criteri di selezione delle espressioni sono i seguenti:

²³Masini (2007), a riguardo, elenca le possibili tipologie di sostantivi N1: quantificatori (*un po' di pane*), classificatori (*un mazzo di rose*), collettivizzatori (*un gruppo di persone*).

- **annotazione su lemma:** la combinazione deve essere formata da tre elementi in sequenza, inclusi nella ricerca nella loro forma lemmatizzata;
- **annotazione PoS:** il primo ed il terzo elemento della combinazione devono essere contraddistinti dal tag specifico **S**, relativo alla categoria dei sostantivi; il secondo deve essere etichettato con tag specifico **EA**, relativo alla categoria delle preposizioni articolate;
- **annotazione sintattica:** la preposizione deve avere il sostantivo quale testa sintattica, mentre il secondo sostantivo deve dipendere dalla preposizione.

L'insieme delle espressioni così raccolto spazia da frequenza 9.958 (*fine dell'anno*) a 288 (*gara del campionato*) ed è incluso in Appendice E. Anche in questo caso lo strumento opera unicamente il test di interrompibilità nella verifica della modificabilità sintagmatica per i motivi esposti nel caso di NPN.

Come si vede dalla Tabella 4.18, le espressioni che esibiscono blocchi variazionali per entrambi gli indici sono di due tipologie.

Espr.	Forma pref.	Freq.	I_{syn}	I_{sub}
massimo di voto	massimo dei voti	343	0	0,009985
forza di ordine	forze dell'ordine	3.285	0,002732	0,009696
ministro di difesa	ministro della difesa	314	0,009464	0,008168
salto in passato	salto nel passato	384	0,002597	0,007416
attacco al suolo	attacco al suolo	503	0,003960	0,007378
età di bronzo	età del bronzo	1190	0,002515	0,006767
signore di guerra	signori della guerra	449	0,008830	0,005821
metro su livello	[NUM] metri sul livello [del mare]	1.002	0,004965	0,004561
lavaggio di cervello	lavaggio del cervello	344	0,008646	0,004218
coppa di mondo	coppa del mondo	448	0,004444	0,003334
età di ferro	età del ferro	681	0,001466	0,003299
campione di mondo	campione del mondo	3.528	0,001133	0,003298
ordine di giorno	ordine del giorno	1.149	0,004333	0,003085
età di oro	età dell'oro	328	0	0,002068
direttore di fotografia	direttore della fotografia	372	0	0,001295
codice di strada	codice della strada	384	0	0,001179
storia di musica	storia della musica	857	0,008102	0,000563
lista di ospite	lista degli ospiti	430	0,002320	0,000430
fisica di particella	fisica delle particelle	303	0,009804	0
lunedì al venerdì	[dal] lunedì al venerdì	518	0,005758	0
sciopero di fame	sciopero della fame	557	0,003578	0
vigile di fuoco	vigile del fuoco	779	0	0
bocca al lupo	[in] bocca al lupo	368	0	0
tecnico di suono	tecnico del suono	316	0	0

Tabella 4.18: Espressioni del pattern NPdN che mostrano valori per gli indici I_{syn} e I_{sub} inferiori a 0,01. Per maggiore chiarezza è riportata, in questo caso, anche la forma maggiormente attestata dell'espressione.

Da un lato compaiono espressioni che si configurano come unità di significato (*tecnico del suono, vigile del fuoco, ordine del giorno*, ecc.) e che mostrano una grande coesione interna. Dall'altro compaiono ancora tre frammenti di locuzioni o costruzioni (*[in] bocca al lupo, [dal] lunedì al venerdì, [NUM] metri sul livello [del mare]*, ecc.) più o meno fisse.

Le Tabelle 4.19 e 4.20 mostrano invece le espressioni che presentano valori di interrompibilità bassi ($I_{syn} < 0,06$), ma sostituibilità medio-alta ($I_{sub} > 0,33$).

Espr.	Forma pref.	Freq.	I_{syn}	I_{sub}
scopo di gioco	scopo del gioco	350	0,059140	0,409963
capitano di squadra	capitano della squadra	389	0,053528	0,395284
presidente di commissione	presidente della commissione	408	0,053398	0,492305
ufficiale di esercito	ufficiale dell'esercito	398	0,052381	0,373792
capo di stato	capo dello stato	472	0,046465	0,656012
ascesa al potere	ascesa al potere	400	0,040767	0,350664
abitante di zona	abitante della zona	401	0,029056	0,345586
fine di campionato	fine del campionato	408	0,026253	0,627120
caduta di regime	caduta del regime	339	0,025862	0,536956
presa di potere	presa del potere	341	0,025714	0,368057
abitante di luogo	abitanti del luogo	346	0,017045	0,741706
resto di gruppo	resto del gruppo	384	0,015385	0,581023
cosa di genere	cosa del genere	692	0,001443	0,340669

Tabella 4.19: Espressioni del pattern NPdN non incluse in locuzioni più ampie che mostrano valori bassi di variazione sintagmatica ($I_{syn} < 0,06$) e valori alti di sostituzione ($I_{sub} > 0,33$).

Espr.	Forma pref.	Freq.	I_{syn}	I_{sub}
parte di persona	[da] parte delle persone	480	0,058824	0,503216
interno di edificio	[all'] interno dell'edificio	416	0,056689	0,541008
interno di chiesa	[all'] interno della chiesa	1125	0,052233	0,455836
occasione di festa	[in] occasione delle feste	441	0,047516	0,394431
conto di fatto	[nel] conto dei fatti	308	0,040498	0,420192
parte di autorità	[da] parte delle autorità	424	0,037471	0,732664
vetta di classifica	[Pd] vetta della classifica	521	0,036969	0,342254
piede di monte	[ai] piedi del monte	334	0,034682	0,368209
parte di pubblico	[da] parte del pubblico	687	0,033755	0,464355
conoscenza di fatto	[a] conoscenza dei fatti	347	0,022535	0,441547
causa di fatto	[a] causa del fatto	491	0,019960	0,348966

Tabella 4.20: Espressioni del pattern NPdN incluse in locuzioni più ampie che mostrano valori bassi di variazione sintagmatica ($I_{syn} < 0,06$) e valori alti di sostituzione ($I_{sub} > 0,33$).

Nella prima delle due è possibile vedere le espressioni che appaiono completamente composizionali, che sembrano includere tuttavia coppie di componenti alle

quali è possibile attribuire un certo grado di familiarità²⁴. Nella seconda sono invece mostrati i frammenti di locuzioni più ampie che ricadono in questo insieme.

La Tabella 4.21, infine, mostra tutte le espressioni caratterizzate da una variazione sintagmatica superiore al 33%.

Espr.	Forma pref.	Freq.	I_{syn}	I_{sub}
zona di città	zone della città	299	0,662528	0,753840
parte di opera	parte delle opere	374	0,445926	0,736719
anno di guerra	[negli] anni della guerra	384	0,419062	0,716404
sponda di fiume	[sulle] sponde del fiume	432	0,451080	0,700239
parte di città	parte della città	832	0,429355	0,670078
periodo di guerra	periodo della guerra	299	0,410256	0,628485
termine di guerra	[al] termine della guerra	562	0,460653	0,555517
fine di secolo	[alla] fine del secolo	1.189	0,828477	0,541397
fine di serie	fine della serie	318	0,331933	0,450027
uscita di album	uscita dell'album	634	0,405811	0,415819
inizio di secolo	inizio del secolo	867	0,830233	0,403614
fine di guerra	fine della guerra	3.234	0,347623	0,371333
corso di guerra	[nel] corso della guerra	843	0,502067	0,365754
inizio di guerra	inizio della guerra	786	0,419926	0,361098
parte di paese	parte dei paesi	602	0,454710	0,352885
riva di fiume	[sulle] rive del fiume	915	0,377127	0,344090
corso di secolo	[nel] corso dei secoli	2.443	0,405596	0,336723
momento di morte	[al] momento della morte	306	0,392857	0,302658
parte di provincia	parte della provincia	1.235	0,338156	0,283329
parte di storia	parte della storia	502	0,333333	0,276760
parte di tempo	parte del tempo	448	0,483276	0,201313
anno di morte	anno della morte	303	0,547085	0,151894
scoppio di guerra	scoppio della guerra	988	0,590551	0,150328
traccia di album	traccia dell'album	324	0,331959	0,107624
battaglia di guerra	battaglie della guerra	363	0,430141	0,096698
inizio di carriera	inizio della carriera	389	0,569690	0,083440
registrazione di album	registrazione dell'album	353	0,477037	0,072281
metà di secolo	metà del secolo	1.744	0,793144	0,036736
titolo di album	titolo dell'album	367	0,404221	0,028770
pubblicazione di album	pubblicazione dell'album	389	0,441092	0,023598
canzone di album	canzone dell'album	359	0,338858	0,005210
decennio di secolo	decennio del secolo	299	0,691753	0,004336
posizione di classifica	posizione della classifica	508	0,427283	0,001932

Tabella 4.21: Espressioni del pattern NPdN che mostrano alti valori di variazione sintagmatica ($I_{syn} > 0,33$).

Come si vede, le espressioni non sembrano, in generale, esibire particolare coesione. La bassa sostituibilità per le espressioni nella parte inferiore della tabella

²⁴Il caso di *capo dello stato* presente in questo insieme è dovuto alla competizione tra la polirematica che identifica sinonimicamente il Presidente della Repubblica Italiana e le occorrenze in cui l'espressione è usata per designare il capo di un qualsiasi stato. L'alta sostituzione del secondo sostantivo con lemmi quali *governo* e *nazione* (che condividono ampiamente lo stesso contesto di *stato*) amplificano la predominanza della seconda accezione dell'espressione.

sembra motivata dalla presenza di costituenti che non hanno generalmente sinonimi (*classifica, decennio, secolo, album*) più che da restrizioni lessicali. Anche qui, infine, compaiono frammenti di locuzioni quali *[al] termine della guerra, [nel] corso della guerra/dei secoli, [al] momento della morte*).

In definitiva il pattern NPdN, escludendo i frammenti di locuzioni più ampie, mostra solo presenza di polirematiche quando entrambi i blocchi variazionali sono attivi e di espressioni fraseologicamente riconoscibili per variazione sintagmatica bloccata ma variazione paradigmatica ammessa.

4.2.5 Analisi sul pattern NPV_{inf}

Le 500 espressioni relative al pattern NPV_{inf} vengono estratte dal corpus in base ai seguenti criteri:

- **annotazione su lemma:** la combinazione deve essere formata da tre elementi in sequenza, inclusi nella ricerca nella loro forma lemmatizzata;
- **annotazione PoS:** il primo elemento della combinazione deve essere contraddistinto dal tag specifico **S**, relativo alla categoria nominale; il secondo deve essere etichettato con tag specifico **E**, relativo alle preposizioni semplici; l'ultimo deve essere contraddistinto dal tag generale **V** relativo alla categoria verbale e dal tag morfologico **f**, che indica il modo infinito.
- **annotazione sintattica:** la preposizione dipende dal nome e il verbo dipende dalla preposizione relativamente ai legami sintattici.

Le espressioni raccolte hanno frequenze comprese tra 1.396 (*modo di fare*) e 63 (*grado di camminare*) e sono mostrate in Appendice F. Per questo pattern il test di variazione sintagmatica cerca unicamente espressioni che risultano interrotte (tra primo e secondo o tra secondo e terzo componente).

Le costruzioni relative alla sequenza NPV_{inf} non sono un modo di formazione standard di sintagmi nominali in italiano, ad eccezione delle espressioni fortemente cristallizzate.

Come si vede dalle tabelle, gran parte delle espressioni che ricadono in questa sequenza sono frammenti di locuzioni del tipo PNPV_{inf}.

Si nota come solo per l'insieme delle espressioni che presenta sia blocchi sintagmatici che paradigmatici (Tabella 4.22), ad esclusione dei frammenti di locuzioni, sia possibile individuare espressioni che rimandano ad unità di significato (*gomma da masticare, associazione per delinquere, macchina da scrivere*, ecc.). Per tutti gli altri insiemi ci troviamo di fronte ad espressioni senza necessità di inquadramento in ambito fraseologico. Nella Tabella 4.23 si vede come la quasi totalità delle espressioni siano in realtà variazioni tematiche sul verbo della locuzione *in grado di*, con la presenza anche di *al fine di, allo scopo di* e *[avere/assumere] il compito di*. Queste espressioni hanno bassissimi gradi di interrompibilità in quanto la locuzione

Espr.	Freq.	I_{syn}	I_{sub}
grado di sopravvivere	124	0,008000	0,099246
ragione d' essere	73	0	0,080800
scopo di visualizzare	485	0	0,026333
associazione a delinquere	199	0,005000	0,014800
procinto di partire	66	0	0,010423
autorizzazione a procedere	138	0,007194	0,001497
patto di avere	964	0,004132	0,000752
associazione per delinquere	135	0	0
macchina da scrivere	104	0	0
gomma da masticare	71	0	0

Tabella 4.22: Espressioni relative al pattern NPV_{inf} con valori per gli indici di modificabilità sintagmatica e paradigmatica inferiori all'1%. In rosso le unità di significato.

non ammette costituenti intervenienti, mentre la variabilità sintagmatica è garantita dalla piena sostituibilità del verbo con i propri sinonimi. *Modo di dire* appare in questo insieme grazie all'alta sostituibilità del verbo con lemmi quali *parlare*, *dichiarare* che, sebbene non ricostruiscano espressioni totalmente sinonime dell'originale, si trovano spesso in contesti similissimi.

La Tabella 4.24, infine, mostra tutte le espressioni con interrompibilità superiore al 20%.

Come si vede, la totalità delle espressioni che compaiono nell'elenco esibiscono alti valori anche di sostituibilità, e si configurano come espressioni libere.

Si può quindi concludere che solo per il blocco di entrambe le variazioni il pattern NPV_{inf} ammetta espressioni polirematiche.

Espr.	Freq.	I_{syn}	I_{sub}
grado di scrivere	64	0	0,938833
grado di coprire	82	0	0,923333
grado di accogliere	84	0	0,914604
grado di provocare	91	0	0,899021
grado di determinare	116	0,008547	0,898731
grado di proteggere	100	0	0,846659
grado di eseguire	345	0,005764	0,842430
grado di individuare	118	0,008403	0,814256
grado di spiegare	132	0,007519	0,805018
grado di distinguere	124	0,008000	0,796673
fine di assicurare	73	0	0,789972
scopo di garantire	105	0,009434	0,785860
grado di sviluppare	236	0	0,783450
grado di valutare	88	0	0,779040
grado di legare	111	0,008929	0,773961
grado di affrontare	215	0,009217	0,773039
grado di comunicare	140	0,007092	0,768021
grado di costruire	107	0,009259	0,767293
grado di spingere	96	0	0,765527
grado di generare	311	0	0,764321
grado di misurare	95	0	0,756162
modo da assicurare	87	0	0,728554
scopo di facilitare	72	0	0,723798
grado di supportare	117	0	0,720574
grado di raccogliere	67	0	0,701622
grado di controllare	310	0,006410	0,692603
grado di trasportare	300	0,006623	0,688663
grado di uccidere	129	0,007692	0,669512
grado di attivare	63	0	0,647476
grado di mostrare	106	0,009346	0,611753
grado di resistere	189	0	0,598402
grado di riconoscere	219	0,009050	0,587014
grado di aiutare	106	0,009346	0,584238
compito di portare	102	0,009709	0,542692
modo di dire	783	0,008861	0,512126
fine di evitare	217	0,009132	0,506155
grado di sparare	102	0,009709	0,413114
grado di contenere	146	0,006803	0,402804
grado di ricostruire	65	0	0,400376
grado di combattere	123	0,008065	0,391498
grado di erogare	579	0,008562	0,375264
grado di influenzare	81	0	0,369154

Tabella 4.23: Espressioni relative al pattern NPV_{inf} con valori inferiori all'1% per quanto riguarda la modificabilità sintagmatica, ma superiori al 33% per la modificabilità paradigmatica.

Espr.	Freq.	I_{syn}	I_{sub}
scelta di fare	83	0,245455	0,749167
tentativo di ottenere	109	0,226950	0,740931
modo di affrontare	89	0,232759	0,683760
decisione di fare	100	0,236641	0,652724
decisione di abbandonare	68	0,218391	0,640165
modo per dire	74	0,260000	0,619741
tentativo di fare	418	0,200765	0,489907
conto di avere	111	0,255034	0,405069
modo per evitare	66	0,232558	0,349966
conto di essere	160	0,234450	0,347920
tentativo di prendere	65	0,207317	0,332385
modo per fare	328	0,220903	0,269108
modo da essere	126	0,397129	0,257498
consapevolezza di essere	74	0,221053	0,251734
paura di essere	75	0,210526	0,210308
paura di fare	68	0,200000	0,120552

Tabella 4.24: Espressioni relative al pattern NPV_{inf} con valori in interrompibilità superiori al 20%.

4.2.6 Analisi sul pattern NN

Le 500 espressioni relative al pattern NN sono state estratte secondo i seguenti criteri

- **annotazione su lemma:** la combinazione deve essere formata da tre elementi in sequenza, inclusi nella ricerca nella loro forma lemmatizzata;
- **annotazione PoS:** il primo ed il secondo elemento della combinazione devono essere contraddistinti dal tag specifico **S**, relativo alla categoria dei nomi comuni;
- **annotazione sintattica:** il secondo sostantivo deve avere il primo quale testa sintattica.

Poiché, in generale, è frequente che il software di tagging grammaticale assegni la categoria nominale a parole straniere o non riconosciute, anche in sequenza, l'estrazione delle espressioni per il pattern NN è viziata da tale anomalia, oltre che da alcuni casi di errata categorizzazione sostantivo-aggettivo. Dopo un'operazione di filtraggio manuale delle combinazioni scorrette, solo 229 delle 500 espressioni estratte sono effettivamente sequenze sensate in italiano. Le espressioni hanno frequenza di occorrenza massima pari 6.542 (*linea guida*) a 173 (*inizio carriera*) e sono riportate in Appendice G.

Il pattern NN non è un pattern appartenente alla sintassi libera dell'italiano, nonostante tale sequenza possa essere riconosciuta in una serie di combinazioni che agiscono in qualità di composti nominali. La non produttività del pattern rispetto a combinazioni standard è testimoniata dal fatto che della lista in input, nessuna delle espressioni risulta avere alti valori per entrambe le modificabilità sintagmatica e paradigmantica.

All'opposto, le combinazioni che esibiscono blocchi per entrambe le modificabilità si configurano come unità di significato (*gas serra*, *scheda madre*, ecc.), come è mostrato dalla Tabella 4.25.

Appaiono in questo gruppo espressioni straniere lessicalizzate e completamente integrate nella lingua italiana (*luna park*, *home page*, *par condicio*, ...). Inoltre, seguendo Baroni *et al.* (2006), è possibile riconoscere diverse tipologie di espressioni composte: (a) *coordinative* (come nei casi di *principe elettore*, *filo conduttore*, *gas serra*, *tenente colonnello*, ecc.), in cui testa e modificatore denotano entità simili o comparabili, unite in una costruzione copulativa; (b) *argomentali* (*raccolta fondi*, *trasporto truppe*, ecc.), per i quali la testa è generalmente un sostantivo deverbale e il modificatore l'argomento del verbo; (c) *specificative*²⁵ (*sala giochi*, *conferenza stampa*, *temperatura ambiente*), in cui la testa non deriva da alcuna deverbizzazione ed ha un significato che viene contestualizzato o specificato dal modificatore.

²⁵ *grounding* in Baroni *et al.* (2006).

Espr.	Freq.	I_{syn}	I_{sub}	Espr.	Freq.	I_{syn}	I_{sub}
sala gioco	219	0,009050	0,008140	domenica sera	471	0,006329	0
inizio stagione	317	0,006270	0,004434	filo conduttore	346	0,005747	0
aiuto regista	235	0	0,004388	par condicio	176	0,005650	0
emergenza rifiuto	201	0,004950	0,003698	forza lavoro	903	0,005507	0
linea guida	6.542	0,000153	0,003431	fan club	191	0,005208	0
centro benessere	361	0,005510	0,001954	mass media	412	0,004831	0
lingua madre	385	0,007732	0,001837	temperatura ambiente	661	0,004518	0
fine stagione	724	0,004127	0,001667	busta paga	223	0,004464	0
campo profugo	362	0	0,000891	effetto serra	479	0,004158	0
banca dato	348	0	0,000265	sabato sera	1.522	0,003274	0
pausa pranzo	180	0	0,000256	pagina web	672	0,002967	0
vano scala	175	0	0,000088	fine settimana	1.279	0,002340	0
spada laser	469	0	0,000084	onda radio	480	0,002079	0
trasporto truppa	343	0,002907	0,000060	fine agosto	501	0,001992	0
martedì sera	200	0,009901	0	stazione radio	753	0,001326	0
fine settembre	201	0,009852	0	soap opera	1.131	0,000883	0
campagna acquisto	207	0,009569	0	anno luce	1.137	0	0
zona retrocessione	213	0,009302	0	file system	703	0	0
diretta tv	217	0,009132	0	punto vendita	686	0	0
serie tv	1.768	0,008969	0	home video	536	0	0
principe elettore	342	0,008696	0	scheda madre	560	0	0
raccolta fondo	236	0,008403	0	gas serra	468	0	0
conferenza stampa	2.370	0,008368	0	computer grafica	294	0	0
sabato pomeriggio	371	0,008021	0	fondo pensione	252	0	0
tenente colonnello	669	0,007418	0	lunedì sera	239	0	0
asilo nido	282	0,007042	0	luna park	201	0	0
miniserie tv	283	0,007018	0	home page	187	0	0

Tabella 4.25: Espressioni relative al pattern NN con valori in interrompibilità e sostituibilità inferiori all'1%.

Esistono in questo insieme anche espressioni relative a costruzioni temporali produttive che non presentano possibilità di modifica sintattica, ma sono sostituibili con una serie di elementi appartenenti ad una stessa classe benché non sinonimi degli originali. Nei casi di *lunedì/martedì/sabato/domenica sera* il primo elemento può essere sostituito da qualsiasi giorno della settimana o anche altri avverbi temporali (*ieri, domani*), mentre il secondo può trasformarsi in una qualsiasi parte del giorno, come testimoniato dall'ulteriore esempio presente in tabella *sabato pomeriggio*. Un altro caso analogo è quello di *fine agosto/settembre* dove il primo elemento può essere scambiato con altri delimitatori temporali (*inizio, metà*), mentre il secondo con i nomi di mesi, anni o anche stagioni e altre entità in classi temporali²⁶. Tali espressioni, in ogni caso, si configurano come locuzioni avverbiali, spesso introdotte da una preposizione (*di sabato pomeriggio, a fine giornata*) e per questo potenzialmente escludibili dalla presente analisi sui sintagmi nominali.

Le poche espressioni presenti nell'insieme caratterizzato dal blocco della sola modificabilità sintagmatica (Tabella 4.26) sono identificabili come formazioni composte non pienamente rigide. Se le espressioni esibiscono generalmente grande sostituibilità del secondo componente, la dispersione sui sinonimi è nulla o molto bassa, ad indicare che c'è interscambio di norma con un solo lemma (*sito internet/web, servizio viaggiatori/passeggeri, classifica marcatori/difensori, corpo vettura/macchina*). Tali espressioni rientrano, in ogni caso, nel gruppo (c) dei composti specificativi visti poc'anzi.

Se invece la sostituibilità è concentrata sulla testa, ciò sta ad indicare che il composto vede il secondo sostantivo in ruolo principalmente aggettivale e per que-

²⁶Cfr. *inizio 2010* (43 occorrenze), *fine estate* (163), *fine giornata* (130).

Espr.	Freq.	I_{syn}	I_{sub}
sito internet	1.410	0,002829	0,628693
figura chiave	176	0,005650	0,625172
servizio viaggiatore	257	0,003876	0,564350
punto chiave	247	0	0,472886
ruolo chiave	330	0	0,379569
traccia audio	198	0	0,358688
classifica marcatore	519	0,003839	0,188281
piano terra	1666	0	0,186722
corpo vettura	545	0	0,167383
episodio pilota	528	0,001890	0,132470

Tabella 4.26: Espressioni relative al pattern NN con valori in interrompibilità inferiori all'1% e valori di sostituibilità maggiori del 18%.

sto utilizzabile con tutti i sinonimi o entità affini della testa (*ruolo/figura chiave, punto/luogo/posizione chiave; episodio/puntata/serie pilota*). Baroni *et al.* (2006) definiscono tali espressioni *attributive*.

Non esistono o sono scarse le espressioni con valori alti di interrompibilità. L'unica espressione che supera la soglia di $I_{syn} > 0,33$ è *inizio secolo* a causa della competizione dei due lemmi tra le espressioni *inizio del secolo* e *di inizio secolo*.

In definitiva, il pattern NN risulta generatore di un insieme di entità inquadrabili nell'ambito di composti sintagmatici, più che di espressioni appartenenti alla sintassi che vedono una qualche forma di irrigidimento e, di conseguenza, acquistano uno statuto fraseologico. Per i composti coordinativi e argomentali la completa rigidità paradigmatica è indice di una composizione ormai avvenuta e stabile e, in particolare per i secondi, non è da escludere che la preferenza lessicale dipenda da legami collocativi esistenti tra il verbo d'origine e l'oggetto. La sostituibilità per i composti specificativi è garantita solo per un nucleo ristretto di espressioni secondo la presente analisi, e riguarda in ogni caso il modificatore, che presenta un singolo sinonimo competitore. Si potrà distinguere, in questo caso, tra composti specificativi *fissi* o *semifissi*. I composti attributivi, infine, giustificano, come già detto, la sostituibilità della testa grazie al ruolo aggettivale assunto dal modificatore.

In Tabella 4.27 è mostrata la categorizzazione proposta per le espressioni del pattern NN.

Modif. sintagmatica	Modif. paradigmatica	Categoria	Esempio
+	+/-	//	//
-	+	Composto attributivo	<i>punto chiave</i>
-	+	Composto specificativo semifisso	<i>servizio passeggeri</i>
-	-	Composto coordinativo	<i>tenente colonnello</i>
-	-	Composto argomentale	<i>raccolta fondi</i>
-	-	Composto specificativo fisso	<i>sala giochi</i>

Tabella 4.27: Categorizzazione delle espressioni nominali relative al pattern NN in relazione al comportamento empirico rispetto alle variazioni sintagmatiche e paradigmatiche.

4.2.7 Analisi sul pattern NCN

Le espressioni relative al pattern NCN sono state estratte secondo i seguenti criteri:

- **annotazione su lemma:** la combinazione deve essere formata da tre elementi in sequenza, inclusi nella ricerca nella loro forma lemmatizzata;
- **annotazione POS:** il primo ed il terzo elemento della combinazione devono essere contraddistinti dal tag specifico **S**, relativo alla categoria dei sostantivi; il secondo deve essere etichettato con tag generale **C**, relativo alla categoria delle congiunzioni;
- **annotazione sintattica:** sia la congiunzione che il sostantivo coordinato devono avere il primo sostantivo come testa sintattica, secondo la convenzione adottata nell'etichettatura sintattica del corpus.

Le espressioni vanno da un massimo di occorrenze pari a 2.631 (*monumenti e luoghi*) a un minimo di 52 (*scultore e architetto*) e sono riportate in Appendice H.

Il pattern NCN in italiano non forma generalmente sintagmi singoli bensì, come è naturale aspettarsi dalla presenza della congiunzione, sintagmi in coordinazione.

L'argomento dei cosiddetti "binomi coordinati" è stato già affrontato in letteratura, in particolare negli studi di Masini (2007, 2008) per l'italiano, a cui si rimanda anche per un approfondimento storico sugli approcci al tema. I suddetti lavori mettono in luce una caratteristica particolarmente interessante in relazione al presente studio, e cioè che il pattern NCN può configurarsi come un particolare tipo di strategia coordinante che risulta in qualche modo anomalo rispetto alla sintassi libera a causa dell'assenza di determinanti, generando combinazioni che possono occupare varie posizioni intermedie tra entità fisse e libere, come verrà confermato tra breve.

Anche in questo caso, come si evince dai dati in Tabella 4.28, il solo insieme delle espressioni che presentano blocchi sia sintagmatici che paradigmatici si configura come il contenitore per le entità cristallizzate (anche in flessione, come mostrato per completezza in tabella) che possono ambire ad uno status fraseologico (*carta e penna, andata e ritorno, botta e risposta, punto e virgola*, ecc., come anche le espressioni istituzionalizzate indicanti titoli o settori di studi, quali *lettere e filosofia* e *economia e commercio*) o alcune espressioni occorrenti in contesti specialistici (*beni e servizi* in testi economici o *guelfi e ghibellini* in testi di storia).

Compaiono anche in questo caso numerosi frammenti di locuzioni più ampie, quali *[a] immagine e somiglianza, [di ogni] ordine e grado, [con] riga e compasso, [a] ferro e fuoco*, ecc. oltre ai casi che vedono la comparsa dei nomi dei mesi: *[tra/in] luglio e/ed agosto, [tra/in] dicembre e gennaio*. Si segnala, inoltre, l'espressione quantificatrice *decine o centinaia* includibile anch'essa tra le locuzioni più ampie a causa della necessità della preposizione *di* di seguire il binomio. *Acqua e sapone* risulta l'unica espressione aggettivale presente.

Espr.	Freq.	I_{syn}	I_{sub}	I_{infl}
luglio e agosto	322	0,082621	0	0
violino e orchestra	66	0,070423	0	0,030303
aspirazione e scarico	54	0,068966	0,002607	0
dicembre e gennaio	54	0,068966	0	0
carta e penna	116	0,064516	0,004581	0,017241
bene e servizio	386	0,063107	0,006485	0,012953
lettera e filosofia	129	0,058394	0	0
punto e virgola	102	0,055556	0	0,058824
decina o centinaio	62	0,046154	0	0,016129
economia e commercio	113	0,042373	0,002144	0,008850
buono e cattivo	69	0,041667	0,004084	0
acqua e sapone	71	0,040541	0,000350	0
giacca e cravatta	145	0,039735	0	0,006897
guelfo e ghibellino	251	0,030888	0	0,007968
riga e compasso	112	0,026087	0,004396	0
andata e ritorno	1.087	0,019838	0	0
stella e striscia	281	0,010563	0	0,003559
asta e bilanciere	98	0,010101	0	0
uso e consumo	216	0,009174	0,002118	0
ordine e grado	126	0,007874	0,007295	0,079365
botta e risposta	148	0,006711	0	0,013514
ferro e fuoco	335	0,005935	0	0
immagine e somiglianza	152	0	0,009793	0
fretta e furia	212	0	0	0

Tabella 4.28: Espressioni relative al pattern NCN con valori di modificazione sintagmatica e sostituibilità inferiori all'1%.

Non esiste per questo insieme un gruppo di espressioni modificabili paradigmaticamente ma non sintagmaticamente.

L'insieme di binomi coordinati che ammette, invece, modifiche di interruzione o inversione, ma non di sostituzione, è costituito da espressioni che esibiscono una preferenzialità nell'abbinamento lessicale (v. Tabella 4.29) quando i componenti non risultino termini specialistici propri di linguaggi di settore (come nei casi di *protone*, *teologia*, ecc.) o nomi temporali (*venerdì*, *sabato*) già di per sé privi di sinonimi. Gli abbinamenti, in tali espressioni, seguono principalmente criteri di opposizione (es. *vita e morte*), associazione (es. *animali e piante*), complementarità (es. *padre e madre*) e consequenzialità (es. *decollo e atterraggio*). Dai dati è visibile che la maggioranza delle espressioni non ha interruzione bloccata, ammettendo quindi modificatori interni o articoli. Poiché l'insieme delle espressioni che costituiscono il campione è stato selezionato in base alla frequenza, ciò significa che esiste una forte compresenza di espressioni a determinante zero (*bare nouns* secondo la terminologia in Masini 2007, ripresa dai *bare binomials* di Lambrecht 1984) e delle stesse espressioni inserite nella sintassi libera grazie ad articoli e modificatori²⁷. Si noti, in ogni caso, come esistono in questo gruppo alcune espressioni (*maschi e femmine*, *pregi e difetti*, *caccia e pesca*, ecc.) che presentano inversione bloccata (la cui modificabilità è quindi dovuta unicamente all'interrompibilità) e che quindi rappresentano un sottoinsieme di maggiore fissità.

L'unico gruppo di espressioni che, nella tabella, è attestata con interruzione bloccata è quella che vede la ripetizione dello stesso lemma in coordinazione con funzione accrescitiva (*ore e/ed ore*, *mesi e mesi*, *giorni e giorni*, *milioni e milioni*, *ragazzi e ragazze*, ecc.). Tali entità esibiscono un indice di inversione pari a 0,5, in quanto il test, nello scambiare i due costituenti coordinati, ritrova sempre lo stesso numero di occorrenze di partenza. La sostituzione risulta, inoltre, bloccata, in quanto la costruzione stessa richiede la duplicazione dello stesso lemma per veicolare la funzione richiesta.

Per quelle espressioni, invece, per cui non sussiste alcun blocco, è verificato che le classi a cui appartengono le due entità risultano omogenee (appartengono a un *frame* condiviso, come già chiarito in Lambrecht 1984 e Masini 2007), ma ciò sembra essere causato unicamente dalla presenza della coordinazione che ha generalmente la funzione di unire entità semanticamente omogenee (cfr. Tabella 4.30). Anche queste espressioni, che vedono alti valori di modificabilità, subiscono una forte competizione con le corrispondenti espressioni comprensive di articoli e modificatori integrate nella sintassi standard. A differenza dell'insieme di Tabella 4.29, invece, esse non esibiscono una preferenzialità lessicale tra i componenti, ammettendo piena sostituzione. È interessante notare, in ogni caso, che compaiono, nel gruppo, espressioni che non ammettono inversione: esse consistono di elementi in palese sequenzialità.

Alla luce di quanto esposto, possiamo schematizzare la categorizzazione per i binomi coordinati come da Tabella 4.31. I *binomi polirematici* comprendono le

²⁷Cfr., ad esempio, *cielo e terra* (73 occorrenze) con *il cielo e la terra* (83).

Espr.	Freq.	I_{syn}^{int}	I_{syn}^{ord}	I_{syn}	I_{sub}
cielo e terra	73	0,598901	0,188889	0,633166	0,009610
libertà e democrazia	57	0,393617	0,424242	0,580882	0,009166
ragazzo e ragazzo	282	0,422131	0,500000	0,633766	0,008978
adulto e bambino	79	0,368000	0,435714	0,575269	0,008341
decollo e atterraggio	79	0,318966	0,347107	0,500000	0,006358
pianta e animale	102	0,250000	0,430168	0,521127	0,006092
donna e uomo	248	0,412322	0,900442	0,906942	0,005889
chitarra e voce	161	0,337449	0,560109	0,640625	0,005752
sceneggiatore e regista	84	0,134021	0,621622	0,642553	0,005692
pregio e difetto	127	0,331579	0,023077	0,341969	0,004942
personaggio e situazione	59	0,233766	0,337079	0,448598	0,004765
animale e pianta	77	0,180851	0,569832	0,607143	0,004671
fauna e flora	85	0,335938	0,826531	0,840525	0,004618
potenza e coppia	69	0,233333	0,310000	0,429752	0,004246
fede e ragione	70	0,125000	0,270833	0,339623	0,004231
uomo o donna	114	0,396825	0,129771	0,446602	0,003748
vita e morte	125	0,723234	0,038462	0,735169	0,003709
fiume e lago	82	0,226415	0,426573	0,508982	0,003334
campionato e coppa	90	0,403974	0,081633	0,433962	0,002973
latino e greco	108	0,211679	0,425532	0,502304	0,002748
giorno e giorno	100	0,264706	0,500000	0,576271	0,002714
film e serie	76	0,315315	0,126437	0,377049	0,002495
riga e colonna	55	0,388889	0,112903	0,432990	0,002042
segno e sintomo	114	0,123077	0,370166	0,421320	0,001849
padre e madre	63	0,892675	0,258824	0,896552	0,001202
anno e anno	274	0,219373	0,500000	0,561600	0,000810
bambino e adulto	61	0,314607	0,564286	0,636905	0,000529
arrivo e partenza	66	0,426087	0,431034	0,600000	0,000366
notte e giorno	59	0,358696	0,856796	0,867416	0
critica e pubblico	109	0,128000	0,726131	0,736715	0
teologia e filosofia	65	0,166667	0,671717	0,691943	0
morto e ferito	226	0,670073	0,008772	0,671033	0
maschio e femmina	368	0,641675	0,013405	0,643411	0
figlio e nipote	130	0,622093	0,103448	0,637883	0
agricoltura e allevamento	65	0,610778	0,084507	0,624277	0
milione e milione	118	0,382199	0,500000	0,618123	0
sabato e domenica	257	0,558419	0	0,558419	0
anno ed anno	56	0,164179	0,500000	0,544715	0
ora e ora	168	0,142857	0,500000	0,538462	0
bianco e nero	61	0,537879	0	0,537879	0
mese e mese	146	0,104294	0,500000	0,527508	0
pagina e pagina	66	0,095890	0,500000	0,525180	0
centinaio e centinaio	125	0,031008	0,500000	0,507874	0
decina e decina	379	0,025707	0,500000	0,506510	0
migliaio e migliaio	182	0,021505	0,500000	0,505435	0
giorno e notte	353	0,459418	0,143204	0,504213	0
ora ed ora	92	0	0,500000	0,500000	0
re e regina	74	0,430769	0,097561	0,463768	0
architetto e ingegnere	52	0,133333	0,395349	0,446809	0
regista e sceneggiatore	173	0,153374	0,378378	0,441296	0
protone e neutrone	66	0,365385	0,153846	0,431034	0
caccia e pesca	57	0,418367	0,033898	0,430000	0
parola e musica	55	0,202899	0,320988	0,421053	0
salita e discesa	54	0,364706	0,114754	0,413043	0
cinema e televisione	166	0,359073	0,097826	0,400722	0
oro e argento	219	0,313480	0,154440	0,389972	0
venerdì e sabato	99	0,388889	0	0,388889	0
filosofia e teologia	133	0,125000	0,328283	0,387097	0
polizia e carabiniere	84	0,115789	0,282051	0,343750	0

Tabella 4.29: Espressioni relative al pattern NCN con valori di modificazione sintagmatica superiori al 33% e valori di sostituibilità inferiori all'1%.

Espr.	Freq.	I_{syn}^{int}	I_{syb}^{ord}	I_{syn}	I_{sub}
amico e familiare	69	0,158537	0,456693	0,507143	0,784170
produzione e commercializzazione	62	0,409524	0	0,409524	0,736572
quotidiano e rivista	102	0,081081	0,301370	0,341935	0,668233
compagno e compagno	108	0,382857	0,500000	0,618375	0,662815
familiare e amico	58	0,236842	0,543307	0,600000	0,661701
infanzia e giovinezza	94	0,408805	0	0,408805	0,639229
cantante e compositore	66	0,195122	0,214286	0,340000	0,631550
scrittore e regista	59	0,202703	0,243590	0,365591	0,631350
fondatore e presidente	111	0,250000	0,271429	0,413793	0,606230
diritto e libertà	54	0,635135	0,239437	0,672727	0,567009
chiesa e monastero	141	0,318841	0,145455	0,389610	0,552641
infanzia e adolescenza	115	0,584838	0	0,584838	0,534380
libro e articolo	66	0,346535	0,410714	0,551020	0,520574
amico e collega	143	0,153846	0,411523	0,468401	0,519873
artista e intellettuale	61	0,140845	0,314607	0,383838	0,508988
progettazione e costruzione	106	0,329114	0,094017	0,372781	0,508848
tempo e luogo	82	0,196078	0,333333	0,426573	0,500794
studente e docente	55	0,191176	0,414894	0,485981	0,487492
strada e piazza	63	0,466102	0,258824	0,550000	0,480223
voce e pianoforte	56	0,263158	0,272727	0,422680	0,478559
cosa e persona	59	0,169014	0,372340	0,443396	0,460910
scrittore e critico	59	0,213333	0,262500	0,385417	0,431243
chiesa e convento	183	0,428125	0,124402	0,471098	0,402652
città e villaggio	86	0,394366	0,353383	0,544974	0,393285
via e piazza	76	0,530864	0,366667	0,631068	0,376667
rivista e giornale	71	0,089744	0,711382	0,719368	0,370045
peso e misura	56	0,631579	0,050847	0,638710	0,365087
partito e movimento	77	0,214286	0,266667	0,388889	0,338254
dimensione e forma	54	0,325000	0,843931	0,854839	0,333425

Tabella 4.30: Espressioni relative al pattern NCN con valori di modificazione sintagmatica e paradigmatica superiori al 33%.

espressioni con blocchi sia sintagmatici che paradigmatici che mostrano un legame fisso e riconoscibile tra i componenti, tale da far considerare l'espressione unitaria. Definiamo *binomi collocativi* quelle espressioni per cui la sostituibilità dei componenti è inibita (e che mostrano quindi particolari legami sul piano lessicale), ma ammettono interrompibilità; a seconda che l'inversione sia inibita o permessa, distinguiamo tra *binomi collocativi fissi* o *liberi*. Le *costruzioni accrescitive*, invece, sono caratterizzate dalla duplicazione di uno stesso lemma (che quindi garantisce una virtuale reversibilità dei componenti) ma non ammettono sostituzione o interruzione. Infine definiamo *binomi liberi* le espressioni che ammettono sia interruzione che sostituzione, riuscendo a distinguere un sottogruppo di *binomi liberi sequenziali* qualora tra i componenti sussista un legame di sequenzialità causale o temporale.

Interr.	Inversione	Modif. paradigmatica	Categoria	Esempio
+	+	+	Binomio libero	<i>amici e familiari</i>
+	-	+	Binomio libero sequenziale	<i>infanzia e adolescenza</i>
+	+	-	Binomio collocativo libero	<i>adulti e bambini</i>
+	-	-	Binomio collocativo fisso	<i>morti e feriti</i>
-	+	-	Costruzione accrescitiva	<i>ore ed ore</i>
+/-	+/-	+	//	//
-	-	-	Binomio polirematico	<i>botta e risposta</i>

Tabella 4.31: Categorizzazione delle espressioni nominali relative al pattern NCN in relazione al comportamento empirico rispetto alle variazioni sintagmatiche e paradigmatiche.

4.2.8 Analisi sul pattern VCV

Le 500 espressioni relative al pattern VCV sono state estratte secondo i seguenti criteri:

- **annotazione su lemma:** la combinazione deve essere formata da tre elementi in sequenza, inclusi nella ricerca nella loro forma lemmatizzata;
- **annotazione PoS:** la sequenza deve essere costituita da un verbo (tag generale **V**), una congiunzione (tag generale **C**) e un verbo (tag generale **V**);
- **annotazione sintattica:** sia la congiunzione che il secondo verbo sono sintatticamente dipendenti dal primo verbo.

Il numero di occorrenze delle espressioni estratte varia da un massimo di 851 (*scrivere e dirigere*) ad un minimo di 19 (*vendere e comprare*). Tutte le espressioni sono riportate in Appendice I.

Il pattern VCV si configura come una sequenza che in italiano può generare in uscita diverse categorie di espressioni²⁸: nominali (*gratta e vinci*) o aggettivali (*mordi*

²⁸A rigore, quindi, tale sequenza non dovrebbe essere inclusa a pieno titolo tra i pattern nominali, rappresentando quella che ha più possibilità di variazione in merito alle categorie in uscita.

e *fuggi*) con funzione verbale cristallizzata²⁹, o binomi coordinati verbali (*andare e venire*) che possono presentare una coniugazione più o meno completa.

In Tabella 4.32 sono mostrate le espressioni che esibiscono entrambi i blocchi di variazione. Come presumibile, è qui che si collocano le espressioni nominali (*gratta e vinci* e *tira e molla*) e aggettivali (*mordi e fuggi* e *usa e getta*) in cui le forme verbali appaiono cristallizzate³⁰.

Sono presenti, altresì, frammenti di locuzioni verbali più ampie, quali *[non/senza] sapere né leggere né scrivere*, *[come CLI] pare e piace* che appaiono fisse all'interno di modi di dire.

Ad eccezione di *salare e pepare* che presenta un blocco anche flessivo in quanto presente il 96% delle volte nella forma *salate e pepate* (strettamente legata al linguaggio culinario), il resto delle espressioni presenti (tutte verbali) è caratterizzato da un numero discreto di occorrenze in diverse flessioni ed include sia espressioni comuni fraseologicamente riconoscibili (*basta e avanza*, *copia ed incolla*), che combinazioni legate ad ambiti settoriali (*validare e standardizzare*, *nominare e revocare* appaiono nel corpus in contesti prettamente tecnici del linguaggio economico e giuridico) e per questo prive di sinonimi.

Espr.	Freq.	I_{syn}	I_{sub}	I_{infl}
nominare e revocare	19	0,050000	0,004709	0,684211
raggiungere o superare	50	0,090909	0,003626	0,740000
salare e pepare	28	0,034483	0,003414	0,035714
usare e gettare	149	0,050955	0,000073	0,040268
provare e riprovare	35	0,027778	0	0,771429
tirare e mollare	94	0	0	0
mordere e fuggire	73	0	0	0
grattare e vincere	24	0	0	0
copiare ed incollare	23	0	0	0
bastare ed avanzare	20	0	0	0,450000
bastare e avanzare	86	0,065217	0	0,290698
validare e standardizzare	19	0,050000	0	0,210526
leggere né scrivere	21	0,045455	0	0,047619
parere e piacere	23	0,041667	0	0
sapere né leggere	26	0,037037	0	0,730769

Tabella 4.32: Espressioni relative al pattern VCV con valori di modificazione sintagmatica inferiori al 10% e paradigmatica inferiori al 5%.

Le espressioni presenti in tutti gli altri gruppi risultano verbali e per i componenti è ancora valido l'assunto che vede i due verbi appartenere generalmente ad un *frame*

²⁹Le cui caratteristiche sono state già oggetto di studio in Masini & Thornton (2007)

³⁰Nissim & Zaninello (2011) notano come alcune di tali forme verbali possano risultare produttive, riportando esempi tratti da corpora del tipo *grattate e vincerete*, *grattati e vinci*, *grattiamo e vinciamo*, ecc. Tali fenomeni, tuttavia, appaiono in proporzioni minime rispetto al numero di occorrenze della forma base e, per le quattro espressioni elencate, solo *usa e getta* presenta 6 occorrenze diverse dalla forma nominale: *usare e gettare* (2), *usati e gettati* (2), *usata e gettata* (1), *usano e gettano* (1).

condiviso grazie al ruolo coordinante della congiunzione.

Le espressioni di Tabella 4.33 non mostrano blocchi paradigmatici, non presentando di conseguenza particolari legami lessicali tra i componenti. La consequenzialità del secondo verbo rispetto al primo suggerisce, inoltre, che l'inibizione delle modifiche sintagmatiche sia principalmente dovuta all'impossibilità di invertire i componenti.

Espr.	Freq.	I_{syn}	I_{sub}
custodire ed esporre	19	0,095238	0,666064
ideare e progettare	22	0,083333	0,665629
travolgere e uccidere	19	0,095238	0,620983
restaurare ed ampliare	42	0,023256	0,563607
ristrutturare ed ampliare	39	0,093023	0,525601
saccheggiare ed incendiare	26	0,071429	0,478122
ideare ed organizzare	21	0,086957	0,390116
conoscere ed apprezzare	123	0,053846	0,331693
riprendere e sviluppare	70	0,054054	0,318978
promuovere ed organizzare	22	0,083333	0,277562
rivedere e correggere	19	0,050000	0,232173
richiedere ed ottenere	35	0,078947	0,218084
torturare ed uccidere	37	0,051282	0,204507

Tabella 4.33: Espressioni relative al pattern VCV con valori di modificazione sintagmatica inferiori al 10% e paradigmatica superiori al 20%.

La Tabella 4.34 mostra invece espressioni per cui la sostituibilità non è permessa e che quindi vedono la comparsa di una certa familiarità di associazione tra i componenti. Sussiste, infatti, l'istituzionalizzazione di coppie preferenziali in particolare in relazioni di opposizione (*salire/scendere*, *vivere/morire*, *accendere/spengere*, *aprire/chiudere*, ecc.), oltre che di semplice associazione (*mangiare/bere*, *suonare/cantare*, *dirigere/sceneggiare*) o di consequenzialità (*condannare e giustiziare*, *smontare e rimontare*, *partecipare e vincere*, ecc.), queste ultime marcate da valori minimi di I_{syn}^{ord} . Il fatto che la modificabilità sintagmatica sia permessa, testimonia come queste espressioni possano andare incontro ad interruzioni (principalmente avverbiali) o, quando non si tratti di espressioni consequenziali, all'inversione dei componenti.

Infine, i casi di ripetizione dello stesso lemma verbale (*lavorare o lavorare*, *svolgere e svolgere*, ecc.) sono relativi a situazioni in cui la voce del verbo appare coniugata in due tempi diversi in funzione di continuità (*lavorava o lavora*, *svolge e svolgerà*). Similmente ai casi analoghi visti per il pattern NCN, il particolare uso di queste ripetizioni al fine di veicolare una specifica funzione impedisce all'espressione di vedere la sostituzione di uno dei due lemmi. Questa volta, tuttavia, l'interruzione (mediante avverbi) è ammessa.

La Tabella 4.35 mostra infine espressioni con alti valori di modificazioni sintagmatiche e paradigmatiche. Ad eccezione delle espressioni sequenziali, che vedono attivo il solo blocco di inversione, non sembrano esistere per questi casi particolari legami fraseologici da mettere in evidenza.

Espr.	Freq.	I_{syn}^{int}	I_{syn}^{ord}	I_{syn}	I_{sub}
aprire e chiudere	186	0,340426	0,065327	0,369492	0,009240
precedere e seguire	75	0,528302	0,038462	0,537037	0,008632
partecipare e vincere	96	0,431953	0	0,431953	0,006615
sposare e avere	123	0,680519	0	0,680519	0,006239
aprire o chiudere	58	0,159420	0,236842	0,333333	0,005617
vivere e morire	97	0,426036	0,030000	0,436047	0,005457
essere o essere	53	0,933417	0,500000	0,937574	0,005073
precedere o seguire	28	0,243243	0,243243	0,391304	0,004941
smontare e rimontare	33	0,352941	0	0,352941	0,003633
salire e scendere	95	0,285714	0,112150	0,344828	0,003541
svolgere e svolgere	23	0,323529	0,500000	0,596491	0,003017
suonare e cantare	111	0,362070	0,527660	0,627517	0,002458
dirigere e sceneggiare	47	0,598291	0,298507	0,656934	0,001190
mangiare e bere	139	0,321951	0,103226	0,371041	0,000562
ammalare e morire	84	0,601896	0	0,601896	0,000360
datare e firmare	25	0,264706	0,870466	0,876238	0
condannare e giustiziare	23	0,747253	0	0,747253	0
recitare e cantare	23	0,233333	0,581818	0,629032	0
amare o odiare	21	0,580000	0,086957	0,596154	0
sposare e divorziare	24	0,529412	0,225806	0,586207	0
lavorare o lavorare	53	0,116667	0,500000	0,530973	0
amare e odiare	33	0,352941	0,232558	0,459016	0
mangiare e dormire	29	0,256410	0,292683	0,431373	0
accendere e spegnere	29	0,355556	0,147059	0,420000	0
importare ed esportare	19	0,321429	0,095238	0,366667	0
ridere o piangere	25	0,166667	0,264706	0,358974	0

Tabella 4.34: Espressioni relative al pattern VCV con valori di modificazione sintagmatica superiori al 33% e paradigmatica inferiori all'1%.

Espr.	Freq.	I_{syn}^{int}	I_{syn}^{ord}	I_{syn}	I_{sub}
disegnare e costruire	33	0,421053	0,029412	0,431034	0,927610
fondare e gestire	26	0,409091	0,037037	0,422222	0,926581
disegnare e realizzare	35	0,527027	0,054054	0,539474	0,912751
scrivere ed eseguire	31	0,483333	0	0,483333	0,911064
comporre e produrre	21	0,603773	0,192308	0,637931	0,908752
vedere e conoscere	19	0,536585	0,344828	0,627451	0,898187
scrivere e illustrare	39	0,610000	0,025000	0,613861	0,885966
scrivere e recitare	22	0,488372	0,083333	0,511111	0,869399
toccare e superare	19	0,406250	0	0,406250	0,868326
cancellare e sostituire	20	0,393939	0	0,393939	0,865842
vedere e ascoltare	23	0,452381	0,342857	0,574074	0,861811
scrivere e suonare	40	0,428571	0,148936	0,480519	0,857268
comporre e pubblicare	20	0,718310	0	0,718310	0,856780
affrontare e vincere	24	0,351351	0	0,351351	0,852749
esistere e essere	34	0,671875	0,100000	0,751825	0,851543
scrivere e fare	38	0,683333	0,095238	0,693548	0,851348
catturare e condannare	43	0,402778	0	0,402778	0,845235
esistere ed essere	63	0,671875	0,10000	0,683417	0,838839
catturare e condurre	25	0,537037	0	0,537037	0,833942
scrivere e leggere	38	0,377049	0,908434	0,913242	0,831064
sviluppare e distribuire	29	0,591549	0,033333	0,597222	0,817146
scrivere e dire	19	0,512821	0,756410	0,806122	0,810591
scrivere e incidere	21	0,543478	0	0,543478	0,799758
catturare e rinchiudere	35	0,375000	0	0,375000	0,791819
creare e produrre	36	0,628866	0,100000	0,643564	0,788151
cantare e danzare	23	0,148148	0,361111	0,425000	0,759033
incidere e pubblicare	21	0,522727	0	0,522727	0,756333
comporre e cantare	47	0,276923	0,216667	0,397436	0,743280
eliminare e sostituire	23	0,520833	0	0,520833	0,740367
restaurare e ampliare	20	0,090909	0,523810	0,545455	0,738253
comporre ed eseguire	38	0,432836	0	0,432836	0,735902
scrivere e realizzare	25	0,537037	0,38462	0,545455	0,719138
attaccare e uccidere	43	0,462500	0	0,462500	0,715980
ideare e scrivere	32	0,085714	0,428571	0,457627	0,706638
saccheggiare e devastare	23	0,041667	0,477272	0,488889	0,704576
creare e sviluppare	29	0,482143	0,093750	0,508475	0,703008
arrestare e rinchiudere	51	0,445652	0	0,445652	0,700497
trovare e uccidere	20	0,583333	0	0,583333	0,681564
creare e pubblicare	23	0,788991	0,041667	0,790909	0,667154
esistere o essere	19	0,472222	0	0,472222	0,651894
scrivere e arrangiare	21	0,475000	0,045454	0,487805	0,647259
arrestare e condurre	49	0,363636	0	0,363636	0,644570
abbattere e sostituire	19	0,387097	0	0,387097	0,642983
arrestare e giustiziare	35	0,385965	0,027778	0,396552	0,629350
ampliare e restaurare	22	0,214286	0,476190	0,541667	0,621115
ampliare e migliorare	22	0,290323	0,120000	0,352941	0,616715
abolire e sostituire	27	0,500000	0	0,500000	0,611235
trasmettere e ricevere	25	0,137931	0,358974	0,418605	0,604206
colpire e uccidere	32	0,360000	0	0,360000	0,596740
sopprimere e sostituire	39	0,380952	0	0,380952	0,587734
crescere e divenire	23	0,520833	0	0,520833	0,585753
arrestare e portare	68	0,381818	0	0,381818	0,580839
studiare e fare	22	0,568627	0,043478	0,576923	0,575741
arrestare e mettere	22	0,450000	0	0,450000	0,574985
distruggere e ricostruire	73	0,638614	0,013514	0,640394	0,574113
produrre e pubblicare	21	0,817391	0	0,817391	0,545301
ampliare e trasformare	20	0,285714	0,166667	0,375000	0,544890
occupare e saccheggiare	19	0,296296	0,208333	0,406250	0,539114
catturare e portare	85	0,433333	0	0,433333	0,524066
invadere e distruggere	21	0,432432	0	0,432432	0,521883
catturare e imprigionare	57	0,337209	0	0,337209	0,520308
catturare e mettere	26	0,333333	0	0,333333	0,504369
demolire e sostituire	31	0,367347	0	0,367347	0,500049

Tabella 4.35: Espressioni relative al pattern VCV con valori di modificazione sintagmatica superiori al 33% e paradigmatica superiori al 50%.

Interr.	Inversione	Modif. parad.	Categoria	Esempio
+	+	+	Binomio verbale libero	<i>saccheggiare e devastare</i>
+	-	+	Binomio verbale libero sequenziale	<i>arrestare e rinchiudere</i>
+	+	-	Binomio verbale collocativo non fisso	<i>ridere o piangere</i>
+	+	-	Costruzione verbale di continuità	<i>lavorava e lavora</i>
+	-	-	Binomio verbale collocativo fisso	<i>vivere e morire</i>
-	-	+	Binomio verbale libero sequenziale	<i>ideare e progettare</i>
-	-	-	Binomio polirematico nominale	<i>gratta e vinci</i>
-	-	-	Binomio polirematico aggettivale	<i>usa e getta</i>
-	-	-	Binomio polirematico verbale	<i>bastare e avanzare</i>

Tabella 4.36: Categorizzazione delle espressioni relative al pattern VCV in relazione al comportamento empirico rispetto alle variazioni sintagmatiche e paradigmatiche.

In definitiva, per il pattern VCV è possibile concludere che all'interno del solo gruppo di espressioni completamente cristallizzate è possibile ritrovare espressioni nominali o aggettivali polirematiche. In tutti i casi, invece, di espressioni che conservano lo status verbale, il blocco paradigmatico sembra essere l'unico elemento che possa identificare preferenzialità lessicali tra i componenti. La non interrompibilità, infatti, non sembra riuscire ad operare una distinzione categoriale, mentre la non inversione è attiva solo in presenza di sequenzialità causale o temporale tra i componenti. In Tabella 4.36 è mostrato uno schema riassuntivo.

4.3 Analisi sul pattern verbale VDN

Come primo approccio anche allo studio variazionale delle espressioni verbali, si è scelto di compiere un'analisi tramite lo strumento computazionale sul pattern VDN, tipicamente associato al sintagma verbale non marcato in italiano.

Come esposto nel paragrafo 3.3.2, per il pattern in questione il test di variazione sintagmatica contempla interruzione, topicalizzazione dell'oggetto, ripresa pronominale anaforica, passivizzazione e relativizzazione. Le 500 espressioni relative al pattern VDN sono state estratte secondo i seguenti criteri:

- **annotazione su lemma:** la combinazione deve essere formata da tre elementi in sequenza, inclusi nella ricerca nella loro forma lemmatizzata;
- **annotazione POS:** la sequenza deve essere costituita da un verbo (tag generale **V**), un articolo determinativo (tag specifico **RD**) e un nome comune (tag specifico **S**);
- **annotazione sintattica:** il sostantivo deve dipendere sintatticamente dal verbo e l'articolo dal sostantivo.

Le espressioni così selezionate variano da un massimo di occorrenze di 10.445 (*prendere il nome*) ad un minimo di 321 (*permettere il passaggio*) e sono incluse in Appendice J. La distribuzione delle espressioni nel piano $I_{syn}I_{sub}$ è mostrata in Figura 4.3.

Come si vede in Tabella 4.37, la fissità paradigmatica rappresenta per le espressioni un forte segno di legame lessicale. Considerando tutte le espressioni con sostituzione inferiore al 10%, si vede che per valori nulli di variazione sintagmatica compare un'espressione completamente cristallizzata (*cessate il fuoco*), che ha anche perso la sua funzione verbale. All'aumentare della variabilità sintagmatica compaiono via via espressioni verbali metaforiche, cristallizzate nei propri legami lessicali ma libere nella coniugazione (*dare i natali*, *aprire il fuoco*, *perdere la vita*, *avere la meglio*, *chiudere i battenti*, *puntare il dito*, ecc.) con una forte tendenza all'unitarietà semantica. Non a caso, infatti, è spesso possibile sostituire intuitivamente a tali espressioni, un verbo singolo che agisce da sinonimo³¹. Se si guarda, invece, ai valori massimi di modificazione sintagmatica, si riscontra la presenza di espressioni che non si configurano come unità di significato ma solo come abbinamenti favoriti tra lessemi (cfr. casi come *segnare il gol* o *deporre l'uovo*) o verbi supporto (*prendere la decisione*).

Esiste, tuttavia, una fascia intermedia, delimitata approssimativamente da valori di I_{syn} che vanno dal 10 al 20% in cui espressioni unitarie e metaforiche da un lato,

³¹Cfr. *dare in natali* con *produrre*; *aprire il fuoco* con *sparare*; *perdere la vita* con *morire*; *avere la meglio* con *vincere*; *chiudere i battenti* con *fallire*; *puntare il dito* con *accusare*.

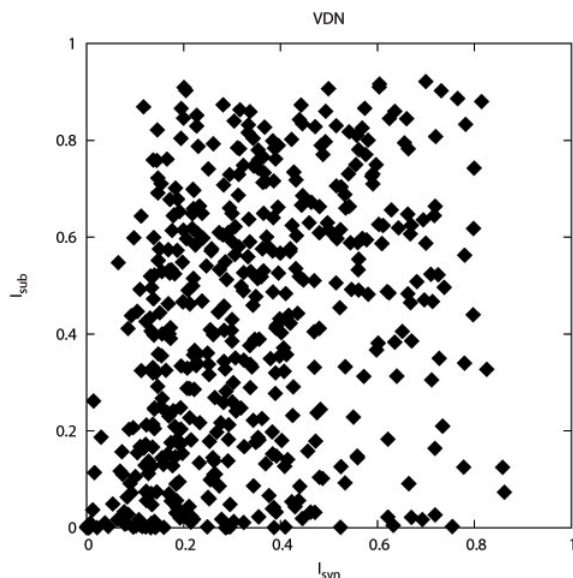


Figura 4.3: Distribuzione delle espressioni in input per il pattern VDN a seconda dei valori empirici ottenuti per gli indici I_{syn} e I_{sub} .

e preferenziali dall'altro, coesistono senza possibilità di discernimento in base anche ad analisi sui singoli test di modificabilità³².

Si segnalano, infine, alcuni casi di espressioni che sembrano generare rumore nella categorizzazione, quali *suonare la chitarra*, *suonare la batteria* e *avere gli occhi*. Per le prime è possibile giustificare la non sostituibilità a causa della mancanza di sinonimi per i termini degli strumenti musicali. La terza espressione, invece, costituisce un frammento della locuzione *avere gli occhi A*, che privilegia in ogni caso il verbo *avere* rispetto a suoi sinonimi come ad esempio *possedere*.

Le espressioni con modificabilità sintagmatica bassa ma sostituibilità medio-alta sono poche, come mostrato dalla Tabella 4.38. Mentre le espressioni con $I_{sub} < 0,2$ possono ancora ambire ad essere incluse nei fenomeni multiparola, le espressioni a più alta sostituibilità appaiono principalmente libere. Ciò conferma ulteriormente che la non modificabilità paradigmatica sia un requisito imprescindibile (e indipendente dalla non modificabilità sintagmatica) per le espressioni verbali al fine di attestare legami fraseologici.

Le espressioni che presentano alti valori di modificazione sintattica e paradigmatica, infine, mostrano comportamenti non riconducibili ad entità multiparola (cfr. Tabella 4.39). Si precisa che le espressioni del pattern VDN in cui è presente il verbo *essere* nel ruolo di copula si concentrano in questo gruppo, mostrando, come

³²A titolo di esempio si confronti *vedere la luce* con $I_{syn}^{int1} = 0,120$, $I_{syn}^{int2} = 0,012$, $I_{syn}^{topic} = 0,002$, $I_{syn}^{anaf} = 0$, $I_{syn}^{pass} = 0,002$, $I_{syn}^{rel} = 0,003$ e *ricevere la visita* con $I_{syn}^{int1} = 0,108$, $I_{syn}^{int2} = 0,031$, $I_{syn}^{topic} = 0$, $I_{syn}^{anaf} = 0$, $I_{syn}^{pass} = 0$, $I_{syn}^{rel} = 0,002$.

Espr.	Freq.	I_{syn}	I_{sub}	Espr.	Freq.	I_{syn}	I_{sub}
pubblicare il album	885	0,718958	0,026635	lasciare il segno	430	0,185606	0,056365
giocare il partita	418	0,670087	0,021623	mettere il mano	850	0,184261	0,090194
segnare il gol	414	0,632327	0,004948	correre il rischio	612	0,184000	0,030624
prestare il voce	380	0,523810	0,000956	avere il occhio	435	0,180791	0,055666
coniare il termine	360	0,461883	0,032266	aprire il occhio	539	0,174579	0,054996
prendere il decisione	847	0,444954	0,019057	dare il dimissione	412	0,169355	0,070254
stare il cosa	399	0,437236	0,056257	dare il caccia	707	0,159334	0,000624
scoppiare il guerra	390	0,429825	0,043382	valere il pena	1.765	0,144864	0,003340
deporre il uovo	390	0,423077	0,053830	fare il conoscenza	628	0,144414	0,074126
vincere il scudetto	392	0,409639	0,001627	tentare il suicidio	389	0,141280	0,000618
togliere il vita	390	0,390625	0,032566	girare il mondo	424	0,139959	0,013269
disputare il campionato	411	0,386567	0,000411	vedere il luce	1.645	0,134666	0,047669
perdere il traccia	448	0,382920	0,039109	ricevere il visita	634	0,133880	0,009326
raggiungere il età	393	0,354680	0,022480	chiudere il occhio	398	0,123348	0,055961
prendere il posto	2.251	0,303312	0,012879	suonare il batteria	343	0,122762	0,075444
fare il spesa	486	0,295652	0,048058	cedere il passo	338	0,122078	0,001436
rischiare il vita	407	0,295652	0,001554	prendere il nome	10.445	0,122048	0,074232
rassegnare il dimissione	614	0,289352	0,000769	prendere il sopravvento	640	0,114799	0,002510
vincere il campionato	2.369	0,286231	0,000320	suonare il chitarra	973	0,110603	0,005502
avere il vantaggio	476	0,273282	0,086084	costare il vita	439	0,107724	0,016922
sposare il figlio	520	0,255014	0,099139	cogliere il occasione	1.120	0,092382	0,072261
portare il squadra	349	0,252677	0,059814	puntare il dito	427	0,089552	0,091662
vincere il serata	743	0,248736	0,001917	riservare il diritto	577	0,085578	0,027837
raggiungere il velocità	393	0,220238	0,017680	trovare il morte	704	0,083333	0,063612
segnare il ritorno	330	0,216152	0,023365	chiudere il battente	353	0,078329	0,008012
prendere il distanza	702	0,214765	0,011565	dare il via	2.108	0,077058	0,032884
vincere il elezione	761	0,211399	0,044183	avere il meglio	1.925	0,063716	0,017166
prevedere il possibilità	392	0,198364	0,028249	perdere il vita	2.034	0,053953	0,049197
passare il notte	582	0,192788	0,020386	aprire il fuoco	709	0,044474	0,002968
catturare il attenzione	335	0,190821	0,094924	dare il natale	691	0,037604	0,009948
raggiungere il semifinale	361	0,186937	0,031605	cessare il fuoco	526	0,001898	0,004675

Tabella 4.37: Espressioni relative al pattern VDN con valori di modificazione paradigmatica inferiori al 10%.

Espr.	Freq.	I_{syn}	I_{sub}
lasciare il gruppo	1.236	0,097151	0,598845
lasciare il scuola	373	0,065163	0,547569
salvare il mondo	354	0,092308	0,438022
perdere il testa	354	0,085271	0,411293
riguardare il composizione	558	0,014134	0,261839
lasciare il band	998	0,096014	0,206564
giungere il momento	583	0,087637	0,203067
perdere il filo	458	0,029661	0,187170
gettare il base	612	0,072727	0,157036
fare il amore	836	0,069042	0,116795
<i>non</i> vedere il ora	926	0,015940	0,113752
prendere il via	1.337	0,074100	0,110082

Tabella 4.38: Espressioni relative al pattern VDN con valori di modificazione sintagmatica inferiori al 10% e valori di modificabilità paradigmatica superiori al 10%.

Espr.	Freq.	I_{syn}	I_{sub}
derivare il nome	702	0,779523	0,339842
andare il cosa	368	0,710008	0,523407
iniziare il carriera	1.393	0,681236	0,508265
ottenere il risultato	351	0,662500	0,483164
continuare il attività	336	0,623740	0,484911
pubblicare il libro	432	0,609050	0,587055
svolgere il attività	415	0,602490	0,626467
iniziare il attività	602	0,601325	0,381201
fare il cosa	520	0,597834	0,367376
portare il nome	1.089	0,590293	0,730660
vincere il gara	521	0,580853	0,482352
scrivere il testo	387	0,546307	0,492255
utilizzare il termine	407	0,532721	0,688259
usare il termine	808	0,514423	0,505692
avere il nome	508	0,492507	0,797401
raggiungere il obiettivo	463	0,479775	0,411928
sorgere il chiesa	325	0,469821	0,331378
svolgere il funzione	690	0,458824	0,629907
vincere il torneo	507	0,438538	0,518820
avere il idea	693	0,432897	0,603402
conquistare il città	435	0,431373	0,535236
usare il nome	362	0,422648	0,802034
attraversare il territorio	344	0,401739	0,402635
continuare il studio	430	0,399441	0,567011
conoscere il nome	364	0,398347	0,431416
riconoscere il diritto	357	0,397976	0,354494

Espr.	Freq.	I_{syn}	I_{sub}
avere il effetto	451	0,394631	0,662226
raccontare il storia	1.632	0,391726	0,417279
affidare il incarico	365	0,389632	0,761618
compiere il studio	484	0,385787	0,716155
raggiungere il livello	355	0,384749	0,640261
scrivere il sceneggiatura	351	0,384211	0,486873
affidare il compito	500	0,383477	0,571629
vedere il successo	373	0,375209	0,580011
lasciare il casa	343	0,364815	0,649509
conquistare il titolo	688	0,364140	0,765766
usare il parola	369	0,363793	0,713169
svolgere il ruolo	590	0,360087	0,518200
dare il possibilità	1.473	0,355643	0,389399
dire il nome	344	0,354597	0,759946
assumere il denominazione	739	0,354021	0,742316
percorrere il strada	379	0,352740	0,514229
seguire il vicenda	416	0,348983	0,477493
leggere il libro	371	0,346831	0,389372
narrare il storia	742	0,344248	0,746908
essere il chiave	392	0,343384	0,804224
proseguire il studio	517	0,343075	0,459361
ridurre il numero	436	0,336377	0,376221
raggiungere il apice	561	0,334520	0,637408
avere il titolo	495	0,333782	0,831979
avere il funzione	920	0,333333	0,815988

Tabella 4.39: Espressioni relative al pattern VDN con valori di modificazione sintagmatica e paradigmatica superiori al 33% senza le espressioni contenenti il verbo essere.

ci si aspetterebbe, alta possibilità di interruzione e sostituzione del predicato. Esse, tuttavia, sono state filtrate dalla Tabella 4.39 al fine di rendere più leggibili i dati³³.

Alla luce di quanto esposto, per il pattern VDN è possibile concludere che la categorizzazione rispetto ai criteri variazionali risulta meno netta rispetto ai pattern nominali e maggiormente orientata all'individuazione di polarizzazioni senza l'esclusione di sovrapposizioni nelle fasce intermedie del continuum variazionale. Per le espressioni che esibiscono blocchi sia sintagmatici che paradigmatici è chiara la possibilità di individuare il polo delle *polirematiche* (sia per espressioni transcategorizzate che per espressioni verbali fisse e richiamanti il concetto di *idiom*). Mantenendo la sostituibilità bloccata, la possibilità delle espressioni di subire modifiche sintattiche o sintagmatiche attesta la comparsa di espressioni con legami lessicali preferenziali o di verbi supporto, che possono confluire sotto l'etichetta delle *collocazioni*. L'assenza di ogni blocco lascia spazio, invece, ad espressioni pienamente integrate nella sintassi libera. La Tabella 4.40 riassume tale schematizzazione.

Modif. sintagmatica	Modif. paradigmatica	Categoria	Esempio
+	+	Espressione libera	<i>scrivere la sceneggiatura</i>
+	-	Collocazione	<i>deporre l'uovo</i>
-	+	Espressione libera	<i>lasciare il gruppo</i>
-	-	Polirematica verbale	<i>dare i natali</i>
-	-	Polirematica nominale	<i>cessate il fuoco</i>

Tabella 4.40: Categorizzazione delle espressioni relative al pattern VDN in relazione al comportamento empirico rispetto alle variazioni sintagmatiche e paradigmatiche.

³³Si ricorda, in ogni caso, che i dati completi per ogni espressione sono riportati in appendice J.

4.4 Conclusioni

Le analisi condotte in questo capitolo hanno permesso di esplorare a fondo il comportamento variazionale di classi di espressioni di diversa natura e raggruppate sotto diversi pattern.

Si è verificato, in generale, che l'assenza di blocchi di modificabilità individua l'insieme delle espressioni libere, evidenziando, all'opposto, come la presenza di almeno un blocco possa individuare espressioni di interesse fraseologico.

Uno dei principali risultati raggiunti è stata la verifica che, indipendentemente dalla categoria in uscita, espressioni appartenenti a diversi pattern mostrano diversi blocchi variazionali, che individuano in modo diverso le categorie in cui è possibile includerle. Ciò amplia la prospettiva spesso ipotizzata in precedenti studi (sull'italiano cfr. Masini 2009, p. 82), in cui, a ciò che qui viene inteso come polirematica, veniva associata una fissità totale, mentre le entità assimilabili alle collocazioni lessicali venivano ricondotte alla sola inibizione della sostituibilità. Se, infatti, è vero che l'assenza di modificabilità contraddistingue in generale le polirematiche (che comprendono unità di significato o espressioni terminologiche), le collocazioni, o in generale gli abbinamenti lessicali preferenziali, si muovono su un terreno più ampio e variabile.

La Tabella 4.41 mostra le combinazioni di modificabilità corrispondenti, per ogni pattern, all'area in cui si concentrano le entità collocative. Dallo schema risulta che il sintagma nominale preferisce instaurare legami preferenziali tra costituenti grazie all'inibizione di variazioni sintagmatiche interne³⁴, mentre nei casi in cui il pattern rappresenti più sintagmi di pari livello (due nominali nel caso di NCN, due verbali per VCV) o un sintagma verbale, essa viene ricondotta a legami lessicali.

Il pattern NPN, il solo che mostra un gruppo di sintagmi nominali collocativi anche nel caso di interruzioni ammesse, vede, in questi casi, la sostituibilità inibita. In questo senso, è interessante notare come il legame fraseologico sia garantito attraverso il ricorso ad almeno una delle inibizioni variazionali rispetto alle espressioni libere. La diversità di comportamento per questo pattern può essere motivata dal fatto che la mancanza del contorno sintattico del determinante può essere già di per sé indice di un certo grado di discostamento dalla sintassi pienamente libera, concedendo all'unità del sintagma di essere spezzata, a patto che la sostituibilità rimanga inibita.

³⁴Dei pattern nominali, le sequenze NA e NPdN risultano le sole combinazioni non marcate e pienamente integrate nella sintassi libera e per questo i migliori oggetti di studio al fine di isolare il nucleo di caratteristiche peculiari delle collocazioni indipendenti da contingenze semantico-sintattiche. Si ricordi, infatti, che AN è una sequenza marcata; NPN una sequenza in cui al nome manca il naturale "contorno sintattico" dell'articolo; NN, NPV_{inf}, VCV sequenze nominali esclusivamente utilizzabili come composti e non formanti collocazioni; NCN una combinazione che mostra interferenza con un particolare tipo di strategia coordinante a determinante zero che si colloca a metà strada tra preferenzialità e sintassi libera Masini (2008).

Pattern	Modif. sintagmatiche	Modif. Paradigmatiche
NA	-	+
AN	//	//
NPN	-	+
	+	-
NPdN	-	+
NPV _{inf}	//	//
NN	//	//
NCN	+	-
VCV	+	-
VDN	+	-

Tabella 4.41: Distribuzione delle entità collocative a seconda delle modificabilità sintagmatiche e paradigmatiche per i pattern analizzati.

All'opposto del comportamento tendenziale dei sintagmi nominali, il pattern verbale VDN, pienamente integrato anch'esso nella sintassi libera italiana, vede concentrarsi le collocazioni verso il polo di modifica paradigmatica bloccata, investendo il livello lessicale del compito di mantenere il legame tra i componenti, e permettendo invece le modifiche sintattiche o sintagmatiche. Quest'ultimo comportamento è riscontrato anche per i pattern NCN e VCV, a indicazione del fatto che i sintagmi distinti di pari livello hanno maggiore facilità di movimento entro i confini della sintassi e investono, quindi, la non sostituibilità del ruolo di garante di un legame in grado di mantenerne la riconoscibilità.

Studio sul pattern NA: il caso della fisica

5.1 Terminologia del linguaggio fisico

In questo capitolo viene esposto un caso di studio sulle espressioni italiane relative al pattern NA e appartenenti al linguaggio tecnico-specialistico della fisica.

La scelta di uno studio su tale linguaggio settoriale risulta interessante a livello linguistico in quanto il lessico fisico, nonostante la sua connotazione tecnica, appare fortemente ricco di parole comuni ed in particolare di espressioni polirematiche formate a partire da queste.

Già Migliorini (1960) nota come la fisica moderna adotti un lessico vicino al linguaggio comune a differenza, ad esempio, della terminologia medica che ricorre a tecnicismi creati attraverso il massiccio utilizzo di radici e affissi non italiani. Tale peculiarità può essere ricondotta alla conservazione delle scelte linguistiche e stilistiche di uno dei padri fondatori della fisica moderna, Galileo Galilei, che ha introdotto una parte cospicua della terminologia:

Il proposito di Galileo di tenere un tono accessibile alle persone colte, anche se non specialiste, ha per corollario il metodo che egli segue quando ha bisogno di termini tecnici: anziché ricorrere al greco o al latino per trarne vocaboli nuovi, preferisce ricorrere a parole usuali, stabilmente adibendole a una nozione specifica. La via scelta da Galileo è ancor oggi, in complesso, quella preferita dai fisici: e una sua influenza in questo campo ci sembra certa (Migliorini, 1960, p. 398).

Casadei (1994) conferma le osservazioni di Migliorini, attestando che su un campione di otto testi di fisica (specialistici o divulgativi), in media il 60% del vocabolario è costituito da parole del vocabolario di base (De Mauro, 1980), ovvero l'insieme dei vocaboli più comunemente noti ai parlanti della lingua, indipendentemente dalla loro collocazione culturale. Come nota l'autrice, in fisica «il tecnicismo è realizzato solo in parte con l'uso di vocaboli estranei all'uso linguistico comune, ed è più spesso costituito dalla riformulazione semantica di vocaboli appartenenti al settore più comune della lingua» (Casadei, 1994, p. 57). Tale meccanismo è accentuato

nella formazione di espressioni polirematiche terminologiche, che vedono spesso la riformulazione semantica e la specializzazione dei lemmi comuni proprio in virtù della loro unione (si pensi a *radice quadrata*, *moto uniforme*, ecc.).

Nell'insieme delle polirematiche del linguaggio fisico, la categoria grammaticale preferita in uscita sembra essere quella nominale. Se si considerano, infatti, tutte le espressioni polirematiche presenti nel GRADIT (De Mauro, 1999-2007) accomunate dall'appartenenza al linguaggio della fisica (etichetta FIS), solo 9 delle 2.677 totali non hanno ruolo nominale (cfr. Tabella 5.1).

Categoria	Lemmi	Esempi
Nominale	2668	<i>corpo nero, funzione d'onda</i>
Aggettivale	4	<i>non lineare, di riferimento, in quadratura</i>
Aggettivo-avverbiale	3	<i>a cascata, in cascata, allo stato libero</i>
Avverbiale	1	<i>all'infinito</i>
Verbale	1	<i>mascherare un suono</i>

Tabella 5.1: Distribuzione in categorie grammaticali delle espressioni polirematiche del linguaggio tecnico-specialistico della fisica estratte dal GRADIT.

All'interno della categoria nominale, inoltre, il pattern più comune per la formazione terminologica è la sequenza NA, come attestato dalla Tabella 5.2, che mostra la distribuzione dei pattern tra le espressioni estratte dal GRADIT con frequenza maggiore di 5.

Tale pattern è stato quindi scelto come soggetto per un primo approccio che studi la possibilità di utilizzare gli indici di variazione sintagmatica e paradigmatica per l'individuazione di terminologia multiparola nel linguaggio fisico, come mostrato nel seguito.

Pattern grammaticale	Occorrenze	Esempio
N A	1.551	<i>bagno termico</i>
N P N	483	<i>forza di gravità</i>
N P N _{pr}	196	<i>ascensore di Einstein</i>
N Pd N	87	<i>spazio delle fasi</i>
N N	65	<i>effetto tunnel</i>
N P N A	54	<i>legge di gravitazione universale</i>
N A A	42	<i>cammino libero medio</i>
N N _{pr}	38	<i>effetto Joule</i>
N Pd N A	19	<i>fisica dello stato solido</i>
N Avv A	12	<i>moto uniformemente accelerato</i>
N A P N	12	<i>asse principale di inerzia</i>
N P D _{ind} N	11	<i>momento di una coppia</i>
N P A N	11	<i>coefficiente di mutua induzione</i>
A N	11	<i>alto vuoto</i>
N P N P N	6	<i>sezione d'urto di cattura</i>
N Pd N P N	6	<i>momento della quantità di moto</i>
N Pd A N	6	<i>legge del minimo sforzo</i>
N A P N _{pr}	6	<i>numero critico di Reynolds</i>
N A Pd N	6	<i>costante universale dei gas</i>

Tabella 5.2: Distribuzione per pattern grammaticale delle espressioni polirematiche del linguaggio tecnico della fisica estratte dal GRADIT. Sono riportati i soli pattern che hanno frequenza maggiore di 5.

5.2 Il corpus

La base empirica necessaria a compiere lo studio sulle espressioni del linguaggio fisico è stata costruita *ad hoc* attraverso la raccolta automatica o manuale di testi e documenti.

Il corpus finale risulta composto da circa 1,8 milioni di *tokens* e raccoglie testi in lingua italiana di varia natura, tutti accomunati dal rappresentare campioni di linguaggio tecnico-specialistico nell'ambito della fisica. Al fine di dare spazio a una maggiore varietà di lingua che non fosse esclusivamente “tecnica” in maniera preponderante, si è scelto di considerare diverse categorie di testi che spaziano da ambiti fortemente specialistici (tesi di laurea in fisica) ad esempi più vicini al linguaggio piano della narrativa (testi divulgativi rivolti ad un pubblico non specialista) i quali, pur privilegiando uno stile più semplice e colloquiale, non mancano della caratteristica componente terminologica a cui è volta la presente analisi, vale a dire le espressioni multiparola. Esse, infatti, come già sottolineato in precedenza, rappresentano una buona parte del lessico terminologico di domini specialistici e, specie in fisica, risultano spesso formate da componenti generalmente presenti in autonomia nel lessico colloquiale.

5.2.1 Costituzione del corpus

I testi raccolti nel corpus sono suddivisi in quattro maggiori categorie che, nell'ordine, possono essere disposte su di un asse che vada da un minimo di tecnicità fino a campioni via via più specialistici:

- testi di divulgazione scientifica;
- voci estratte da Wikipedia;
- manuali specialistici universitari;
- tesi di laurea triennali, specialistiche e magistrali.

In Tabella 5.3 sono esposti i dati quantitativi relativi alla suddivisione del corpus in termini di *tokens* e di percentuali totali delle singole categorie; queste ultime sono discusse in dettaglio nel seguito.

La scelta di non includere la letteratura scientifica è derivata dal fatto che la quasi totalità di articoli e ricerche di settore (come anche gran parte dei manuali e testi normalmente in uso nelle università) è pubblicata in lingua inglese che, in questo ambito come in molti altri, è divenuta al giorno d'oggi una vera *lingua franca*.

5.2.1.1 Testi divulgativi

La sezione divulgativa del corpus conta poco più di 100.000 *tokens* e raccoglie testi rivolti ad un pubblico prettamente non specialista. Essi sono caratterizzati, quindi, oltre che da un lessico meno tecnico, anche da un periodare maggiormente

Categoria	N. di <i>tokens</i>	Percentuale
Testi divulgativi	109.209	6,2%
Voci di Wikipedia	608.349	34,5%
Manuali	364.387	20,7%
Tesi di laurea	679.195	38,6%
Totale	1.761.140	100%

Tabella 5.3: Dati quantitativi sulla suddivisione in categorie del corpus.

vicino alla narrazione. Il seguente estratto risulta rappresentativo del campione (in corsivo le espressioni polirematiche):

“[...] Eravate entrati nel *campo magnetico* della Terra, là dove le *linee di forza* cominciano ad infittirsi. Tu sai che una *carica elettrica*, quando è soggetta a un *campo magnetico* devia la traiettoria, in versi opposti se le cariche sono diverse.”

“Ah già... ma questo venni a conoscerlo in seguito. Allora fui trascinato in uno stravagante percorso elicoidale. La qual cosa mi permise di osservare tranquillamente il vostro pianeta Terra... Mio Dio, quale spettacolo!”

(da Manduchi - *Vita e Morte di un elettrone*, 2011)

Se il testo sopra citato può chiaramente apparire come un esempio di divulgazione diretto soprattutto a fruitori giovani, in questa categoria sono stati inclusi anche testi rivolti ad un pubblico adulto, come mostra il seguente esempio:

Con la *relatività generale* e con la *meccanica quantistica* il peso dell'apparato matematico nella fisica diventa sempre più rilevante, meno intuitivo e tecnicamente più sofisticato (si richiedono strumenti non elementari come la *geometria differenziale*, che a partire dalla relatività si svilupperà enormemente, e l'*analisi funzionale* negli *spazi di Hilbert* per la *meccanica quantistica*): la fisica si distanzia dall'esperienza diretta e dal *senso comune*.

(da Strumia - *La teoria della relatività*, 2010)

5.2.1.2 Voci di Wikipedia

Per questa sezione sono state considerate 535 pagine web dell'enciclopedia libera Wikipedia. L'idea di ricorrere a questa risorsa per costruire una parte sostanziale del corpus (circa 600.000 *tokens*) è motivata dalla grande disponibilità e fruibilità di testo in lingua italiana nel settore della fisica. Il carattere enciclopedico della risorsa

è, del resto, un'ottima caratteristica per il campione, che viene a collocarsi in una posizione più alta nella scala di tecnicità rispetto alla divulgazione, rasentando (e in alcuni casi eguagliando) il registro dei manuali specialistici. A titolo d'esempio si riporta di seguito un estratto:

Il teorema di *equipartizione dell'energia* permette di valutare l'entità dell'*energia interna* di un *sistema termodinamico* sulla base di una trattazione classica, non considerando dunque la *quantizzazione dell'energia*: essa è fondata sulla *meccanica statistica classica*, cioè la descrizione newtoniana, o descrizioni più generali, come la formulazione hamiltoniana, con particolare riferimento alle ipotesi della *teoria cinetica dei gas*.

(dalla voce “Teorema di equipartizione dell'energia”)

Le pagine web sono state estratte automaticamente¹ attingendo dalla lista dei link alle voci del Portale Fisica dell'enciclopedia. L'intero contenuto di ogni pagina è stato salvato come testo in singoli file che, come precisato nel seguito, hanno subito una fase di selezione e pulitura.

5.2.1.3 Manuali specialistici

Questa sezione comprende più di 350.000 *tokens* ed è costituita da manuali universitari e dispense normalmente in uso nei corsi di laurea triennali e magistrali in fisica. Il testo, benché con primarie mire esplicative, risulta tecnico e rivolto ad un pubblico con conoscenze di base o consolidate della materia. Di seguito si riporta un estratto del campione:

Per ottenere una lunga durata del criogeno si dovrà evidentemente ridurre l'*ingresso termico* nelle sue tre forme: convettivo, conduttivo e radiativo. Il sistema più semplice per fare ciò è quello di usare un *vaso di Dewar*, o dewar: un contenitore in vetro a doppia parete, mostrato in fig.4.4A, con le pareti argentate (simile ad un thermos).

(da De Bernardis, *Dispense di Laboratorio di Astrofisica*, 2005).

5.2.1.4 Tesi di laurea

La sezione più corposa del corpus (circa 680.000 *tokens*) è formata da 14 tesi di laurea triennale e 21 tesi di laurea specialistica o magistrale discusse tra il 1998 e il 2013 presso il Dipartimento di Fisica dell'Università di Roma “Sapienza”. Tale campione rappresenta la varietà più altamente tecnica presente nel corpus, essendo i testi rivolti ad un pubblico specialista altamente preparato in materia, come mostra il seguente estratto:

¹Tutte le pagine risalgono all'ultima versione disponibile al giorno del salvataggio, 29 maggio 2012

È possibile utilizzare la quantità l_a come *righello standard*, sempre di tipo statistico, essendo essa espressa in termini di quantità costanti e ben definite. Per questo si è soliti farvi riferimento con il nome di *scala acustica*. La *scala acustica*, perciò, è un buon parametro per valutare l'evoluzione del *fattore di scala*. Infatti $r_s(z_r)$ non dipende dalla densità di *energia oscura* né da quella di *curvatura* [...], ma la *distanza di diametro angolare* alla ricombinazione dipende, invece, dall'intera *storia di espansione cosmica*, da oggi a z_r .

(da Salvatelli - *Vincoli cosmologici su modelli di gravità modificata*, Tesi di laurea magistrale in Fisica Teorica, 2012)

5.2.2 Trattamento

5.2.2.1 Processamento pre-tagging

Al completamento della raccolta dei testi componenti il futuro corpus è stato necessario un trattamento di “pulitura” su più fronti, la cui prima fase è stata l'uniformazione del materiale al formato grezzo di testo in codifica UTF-8. Per i testi di divulgazione, i manuali e le tesi di laurea, si è effettuata una conversione dal formato .pdf al formato .txt, con conseguente nascita di errori di ricodifica. Gli interventi maggiori sono quindi stati la sostituzione di alcuni caratteri anomali (ad es. apostrofi, combinazioni di caratteri non UTF-8 in luogo delle lettere accentate, caratteri ASCII) con gli omologhi corretti in UTF-8.

Per le tesi di laurea, inoltre, si è scelto di eliminare parti dei frontespizi (nome di università, facoltà, dipartimenti, ecc.) in quanto elementi di questo tipo avrebbero inficiato una corretta estrazione di espressioni polirematiche a causa della loro alta frequenza e ricorsività. Si sono inoltre eliminate le sezioni bibliografiche finali, per non introdurre nel corpus grandi quantità di nomi propri e titoli stranieri. Per i testi di Wikipedia si è proceduto all'eliminazione dei segmenti ricorrenti presenti in tutte le pagine, quali l'elenco delle lingue, gli *header* e altre informazioni ricorrenti in colonna per i motivi sopra citati.

5.2.2.2 Tagging

Al fine di poter essere utilizzato agli scopi del presente lavoro, il corpus sopra descritto ha subito un trattamento di *part-of-speech tagging* che assegnasse ad ogni *token* una categoria grammaticale e un lemma di riferimento. A tale scopo è stato utilizzato il software TreeTagger (Schmid, 1994), il cui *tagset* è incluso in Appendice K.

Il processo di *tagging* ha condotto alla completa annotazione con lemma e categoria grammaticale di 2.044.956 entità (inclusa la punteggiatura). Si è scelto, inoltre, che il software etichettasse i *tokens* non riconosciuti (non riconducibili, cioè, ad alcun lemma di riferimento presente nel suo dizionario) con la marca *unknown*.

5.2.2.3 Processamento post-tagging

Al fine di migliorare la qualità dell'output del software di *tagging* e gestire, quindi, una risorsa il più possibile accurata, è possibile intraprendere diversi processi semi-automatici o manuali finalizzati alla correzione del materiale processato.

L'importanza di una correzione post-*tagging* risiede nel fatto che un corpus pulito e di grande accuratezza è fondamentale per la qualità delle analisi che in seguito saranno svolte sulla risorsa. In generale un *tagger* può incontrare problemi nel processamento di testi di settore, come i campioni del corpus in esame, a causa del lessico specialistico che può non essere presente nel dizionario interno al software, o per costruzioni grammaticali e sintattiche tipiche dell'ambito su cui il software non ha effettuato uno specifico *training*. Nel seguito è mostrato, tuttavia, che grazie alla sola correzione manuale effettuata a valle del processo di *tagging* è possibile comunque raggiungere un alto livello di accuratezza, evitando di effettuare un *training* specifico per il *tagger* su un *gold standard* di testo specialistico specificamente costruito.

La prima parte del processamento manuale ha riguardato il più ampio recupero possibile dei *tokens* non riconosciuti ed etichettati quindi come *unknowns*, che nel file risultante dal processo di *tagging* risultano pari a 303.515. Si è scelto di estrarre automaticamente tali forme e ordinarle per occorrenze, in modo da individuare i *types* non riconosciuti più frequenti. Si è scelto, quindi, di effettuare un recupero manuale su tutti gli *unknowns* con frequenza $f \geq 8$, per un totale di 227.754 forme (75% degli *unknowns*). Si è proceduto ad uno spoglio manuale per salvare gli *unknowns* validi (*token* chiaramente riconoscibili come parole di senso compiuto ma non individuate dal *tagger*, come ad es. nomi propri, errori di battitura, errori di codifica, lemmi non presenti in dizionario) dalle forme da scartare (parti di formule, pezzi di codice informatico, frammenti di parole o parole interrotte e non riconducibili a lemmi). Nel primo gruppo si sono quindi separate le forme chiaramente riconducibili ad un solo lemma e categoria grammaticale, da quelle ambigue riconducibili, senza contesto, a due o più possibili lemmi. Per le prime si è proceduto ad assegnare a tutte le occorrenze della stessa forma il lemma e la categoria corretta con sostituzioni automatiche. È stato questo il caso di sostantivi o nomi propri non presenti nel dizionario di partenza interno al software, come ad es. *muoni* (576 occorrenze, sostantivo maschile, lemma: *muone*), *autovalore* (132 occorrenze, sostantivo maschile, lemma: *autovalore*), *quantistica* (829 occorrenze, aggettivo femminile, lemma: *quantistico*), *Dirac* (249 occorrenze, nome proprio, lemma: *Dirac*). In questo modo si sono recuperati 1.219 *types* per 52.231 occorrenze totali.

In casi, invece, del tipo *hamiltoniana* (245 occorrenze), non è possibile capire, senza contesto, se la forma si riferisca al sostantivo femminile *hamiltoniana*, o sia aggettivo femminile del lemma *hamiltoniano*. In questi casi si procede ad un laborioso, ma inevitabile, controllo manuale in contesto, assegnando ad ogni forma il lemma e la categoria corretta, riuscendo a recuperare 102 *types*, per un totale di 4.268 occorrenze.

Un'altra importante fase di recupero riguarda quelle forme cui il *tagger* non è stato in grado di assegnare un lemma univoco e che sono etichettate, quindi, con due o più lemmi possibili. Questo genere di situazione riguarda 305 *types* per un totale di 11.490 occorrenze. È questo il caso di *stato* (1164 occorrenze ambigue, in cui il lemma di riferimento può essere il verbo *stare*, il verbo *essere* o il sostantivo *stato*) o *fisica* (868 occorrenze contese tra il sostantivo *fisica* e l'aggettivo *fisico*). Anche in questi casi si effettua un controllo manuale al fine di assegnare alla forma il lemma corretto e si procede fino all'esaurimento di tutte le forme in lista, indipendentemente dalla frequenza di occorrenza.

Un'ulteriore fase di correzione ha riguardato errori macroscopici venuti alla luce nelle suddette fasi di spoglio manuale. Un errore comune del *tagger* è stato quello di attribuire a forme come *questo*, *quello*, *tale*, *alcuno* la categoria di pronomi anche in casi in cui il loro ruolo era quello di aggettivo. Per ovviare a ciò si è automaticamente corretta la categoria in A in tutti i casi in cui il lemma seguente fosse un sostantivo. Similmente si sono corretti casi di congiunzioni come *quindi*, *cioè* o *infatti* etichettate sempre come avverbi per un evidente errore del lessico presente nel *tagger*.

5.2.2.4 Recupero degli unknowns

Il numero totale di *unknowns* presenti nel corpus dopo il processo di *tagging* è pari a 303.515. Al termine di tutte le fasi di recupero è stato ricondotto a lemma un totale di 55.181 *tokens*, il che lascia un numero di *unknowns* ancora presenti pari a 248.334. Di questi, 172.573 hanno una frequenza $f \geq 8$, che li identifica come elementi da scartare, in quanto sono già passati attraverso lo spoglio manuale. I restanti 75.761 comprendono invece sia materiale potenzialmente recuperabile che elementi da non lemmatizzare, ma tali *unknowns* non vengono sottoposti ad alcun recupero manuale per evidenti motivi di impegno e tempo. Tuttavia, è utile provare ad effettuare una stima di quanto materiale valido viene perso conseguentemente a tale scelta e si decide, perciò, di estrarre randomicamente per tre volte un campione di 400 *tokens* tra gli *unknowns* con frequenza $f < 8$. Si procede, dunque, ad un controllo manuale di quanti *tokens* siano recuperabili² ottenendo, sui tre campioni, una media di circa il 18% del totale. In virtù dell'estrazione randomica nella creazione dei tre set di elementi è lecito trasporre tale percentuale al numero totale di *unknowns*, il che fornisce una stima di circa 13.673 *tokens* validi ma non ricondotti a lemma. Conseguentemente il numero totale di *unknowns* da scartare è di circa 234.697, pari a circa l'11% del totale degli elementi lemmatizzati.

5.2.2.5 Accuratezza del tagging

Al fine di valutare i miglioramenti conseguenti ai diversi interventi post-*tagging* sul corpus si è scelto di estrarre in maniera randomica 100 frasi dall'output del

²Con "recuperabile" si intende ogni forma che sarebbe stata ricondotta a lemma se sottoposta al procedimento di recupero effettuato per i *tokens* con frequenza maggiore a quella di soglia.

software di *tagging* (che non comprende, quindi, alcuna correzione manuale) e dal file contenente il corpus finale, che considera invece tutti i processi correttivi. In questo modo si hanno a disposizione due campioni rispettivamente di $n_i = 2.998$ e $n_f = 3.417$ *tokens*, che, in virtù dell'aleatorietà di costruzione, possono fornire informazioni sull'intero corpus. Si procede quindi ad un controllo manuale sulle occorrenze di ciascun campione.

Il primo file, generato al termine del processo di *tagging*, comprende 2.513 *tokens* correttamente categorizzati e lemmatizzati (83,8%), 36 *tokens* etichettati con lemma ambiguo (1,2%), 117 *tokens* a cui è stata attribuita un'errata categoria grammaticale o lemma (3,9%) e un numero totale di *unknowns* pari a 332 (11,1%), di cui 113 recuperabili (cioè forme riconducibili ad un corretto lemma) e 219 da scartare (relativi a parti di formule, abbreviazioni, parole straniere, ecc.).

Con questi dati a disposizione è possibile calcolare la *precisione* iniziale p_i , vale a dire una quantità che valuti quanto bene il *tagger* ha individuato categoria grammaticale e lemma per il materiale riconosciuto. Detto tp il numero totale di *tokens* correttamente lemmatizzati (*true positives*) e ip il numero di *tokens* cui è stato attribuito una categoria errata o un lemma errato o ambiguo (*incorrect positives*), la precisione è data dalla seguente formula:

$$p_i = \frac{tp}{tp + ip} = \frac{2.513}{2.666} = 94,3\% \quad (5.1)$$

È possibile, inoltre, calcolare il *recall* iniziale r_i , vale a dire una misura di ciò che il *tagger* ha riconosciuto e lemmatizzato rispetto al totale del materiale lemmatizzabile. Detto fn il numero totale dei *tokens* da lemmatizzare ma etichettati come *unknown* (*false negatives*), il *recall* risulta:

$$r_i = \frac{tp + ip}{tp + ip + fn} = \frac{2.666}{2.779} = 95,9\% \quad (5.2)$$

Infine è possibile avere una stima dell'*accuratezza* iniziale totale a_i raggiunta dal *tagger*. Detto tn il numero di *unknowns* effettivamente da scartare (*true negatives*), l'*accuratezza* iniziale è la seguente:

$$a_i = \frac{tp + tn}{tp + ip + fn + tn} = \frac{tp + tn}{n_i} = \frac{2.732}{2.998} = 91,1\% \quad (5.3)$$

Considerando, invece, il secondo file contenente le frasi estratte randomicamente dal corpus a valle di tutte le fasi di correzione, si ha un numero totale di *tokens* correttamente categorizzati e lemmatizzati pari a 2.909 (85,1%), 112 *tokens* a cui è stato attribuito un lemma o una categoria errati (3,3%) e 396 *tokens* etichettati come *unknowns* (11,6%), di cui 36 recuperabili e 360 da scartare. Analogamente a quanto fatto per il campione iniziale, è possibile calcolare precisione, *recall* e accuratezza finali (rispettivamente p_f , r_f e a_f) con l'utilizzo delle stesse formule viste sopra. I risultati sono i seguenti:

$$p_f = \frac{2.909}{3.021} = 96,3\% \quad (5.4)$$

$$r_f = \frac{3.021}{3.057} = 98,8\% \quad (5.5)$$

$$a_f = \frac{3.269}{3.417} = 95,7\% \quad (5.6)$$

Come si vede dal confronto dei risultati, dopo le fasi di correzione semi-automatica e manuale si ha un notevole incremento di tutti e tre gli indici, il che comporta un sensibile miglioramento della risorsa e, di conseguenza, un aumento dell'affidabilità delle analisi che in seguito verranno effettuate su di essa.

5.3 Caratteristiche delle polirematiche fisiche NA

Grazie al corpus descritto, è possibile indagare le proprietà variazionali delle espressioni fisiche corrispondenti al pattern NA grazie allo strumento computazionale e agli indici di modificabilità relative al pattern (interruzione, ordine dei costituenti e sostituzione con sinonimi).

Se, da un lato, i due test di modificabilità sintattica non pongono problemi di applicazione al linguaggio settoriale in questione rispetto a quanto fatto per il linguaggio generale nel precedente capitolo, la peculiarità del lessico fisico di essere costituito in larga parte da vocabolario comune rende possibile anche l'applicazione del test di sostituibilità, cosa non ovvia, in generale, per qualsiasi linguaggio specialistico. Un'alta proporzione di termini tecnici, infatti, comporta in generale una bassa o nulla presenza di sinonimi per questi ultimi, rendendo, di fatto, il test inefficace.

Alla luce delle analisi esposte nel precedente capitolo, è ragionevole ritenere che la terminologia fisica multiparola sia costituita da espressioni che presentino blocchi sia sintagmatici che paradigmatici.

Per testare tale ipotesi si è studiato il comportamento variazionale delle polirematiche fisiche del GRADIT relative al pattern NA presenti nel corpus, risultate nel numero di 595 (su 1.551 totali censite nel GRADIT). In questo senso si considera la lista del GRADIT un *gold standard* di riferimento, che includa espressioni terminologiche con uno status categoriale definito.

Dati i due soli livelli di annotazione presenti per il corpus (lemma e PoS), non è possibile effettuare le query dello strumento computazionale includendo anche l'informazione sintattica. Per garantire affidabilità al test di interruzione si è scelto quindi di limitare ad una sola parola il materiale linguistico interveniente tra i componenti.

La Tabella 5.4 mostra i risultati sui blocchi variazionali per le espressioni della lista estratta dal GRADIT.

Come si vede, i risultati hanno mostrato che il 73% delle espressioni non viene mai interrotto, il 93% non appare mai in ordine inverso e il 64% non ha attestazioni di casi in cui uno dei componenti viene sostituito da un sinonimo. Queste percentuali, come è ovvio, salgono all'aumentare della soglia limite considerata per ognuno degli indici. Esse potrebbero, però essere viziate dalle basse frequenze di molte delle

Espr. totali (595)				Espr. con $f > 30$ (104)			
	$I = 0$	$I < 0,1$	$I < 0,5$		$I = 0$	$I < 0,1$	$I < 0,5$
int	73%	82%	95%	int	54%	90%	100%
ord	93%	97%	99%	ord	85%	100%	100%
sub	64%	70%	85%	sub	49%	73%	95%

Tabella 5.4: Percentuali di espressioni fisiche della lista GRADIT per diversi valori per gli indici I di variazione empirica ricavati dal corpus di fisica (int = interruzione; ord = inversione; sub = sostituzione). A sinistra sono riportati i dati per l'intero insieme di espressioni. A destra i dati delle sole espressioni con numero di occorrenze maggiori di 30.

espressioni. Solo 104, delle 595 totali, infatti, hanno frequenza maggiore di 30. Se si considerano, tuttavia, le sole espressioni con frequenza di occorrenza superiore a 30, le proporzioni rimangono coerenti (cfr. la parte destra di Tabella 5.4).

L'evidenza empirica, quindi, fornisce alte percentuali di blocco variazionale per le espressioni nominali del pattern NA, confermando che il prototipo di tali espressioni terminologiche tende verso il polo di non modificabilità già identificato nelle analisi sul pattern NA del linguaggio generale. Nonostante le consistenti percentuali di blocco sui tre indici, tuttavia, i dati suggerirebbero che la fissità sintagmatica (ed in particolare la fissità dell'ordine dei costituenti) sia maggiormente caratteristica di queste espressioni rispetto al blocco paradigmatico.

È necessario, tuttavia, fare alcune precisazioni.

Il blocco dell'inversione rispetto all'ordine dei costituenti risulta la caratteristica più stabile poiché, oltre ad essere un tratto tipico delle espressioni identificanti unità di significato, esso si può identificare all'interno della sintassi e della semantica italiana come una strategia per garantire obiettività al significato del modificatore³, caratteristica, questa, imprescindibile nella formazione di terminologia specialistica.

Le espressioni dell'insieme che risultano, invece, maggiormente interrompibili, vedono nella maggior parte dei casi l'intervento tra i componenti di un avverbio che genera varianti di similarità o opposizione rispetto all'espressione originale, come nei casi di *campo (non) uniforme*, *campo (quasi) statico*, *trasformazione (non) isoterma*, ecc. In altri casi è l'intervento della copula a provocare l'interruzione (cfr. *grandezza scalare* vs. *questa grandezza è scalare*). In questi casi, le varianti interrotte testimoniano un certo grado di indipendenza dei costituenti rispetto alle espressioni completamente fisse⁴ che le configura come entità tendenti alle collocazioni, pur rientrando nell'ambito di interesse della terminologia.

³Come già discusso nel par. 4.2.2, l'anteposizione del modificatore rispetto alla testa porta con sé varianti stilistiche, marcate o giudizi soggettivi del parlante, rispetto alla costruzione con modificatore post-nominale. Non a caso le polirematiche nominali fisiche censite nel GRADIT e relative al pattern AN rappresentano meno dello 0,5% del totale.

⁴cfr., ad esempio, *equazione caratteristica* vs. **l'equazione è caratteristica*.

La variazione paradigmatica, infine, viene garantita dalla possibilità, per alcune espressioni, di sostituire un componente con un sinonimo stretto (*gas ideale/perfetto*, *campo centrale/radiale*) o con un lemma che nel linguaggio generale rappresenterebbe un sinonimo, ma che in ambito specialistico, a causa della generale marcata mono-referenzialità, si riferisce a concetti affini ma non identici (*forza/energia elettrica*, *meccanica/dinamica newtoniana*).

5.4 Indice di prototipicità

I dati scaturiti dallo studio delle variazioni delle espressioni fisiche estratte dalla lista del GRADIT mostra che il 70% delle espressioni con frequenza maggiore di 30 presenta valori di variabilità inferiori al 10% per tutti e tre gli indici. Tale tendenza, in accordo con quanto rilevato anche nell'analisi del pattern NA del linguaggio generale, mostra come il prototipo di polirematica terminologica relativa al pattern NA sia localizzato nel polo di totale fissità sia sintagmatica che paradigmatica.

In quest'ottica, se si vuole procedere ad un tentativo di identificazione automatica delle espressioni di questo tipo che sfrutti le informazioni variazionali fornite dallo strumento computazionale, è possibile considerare una funzione che riassume in un solo indice i dati ottenuti dai valori di I_{syn}^{int} , I_{syn}^{ord} e I_{sub} e che misuri l'aderenza di una data espressione al prototipo ipotizzato, che preveda totale fissità.

La funzione proposta in questo lavoro è l'Indice di Prototipicità (IP), i cui valori sono forniti dalla seguente formula:

$$IP = \frac{n_{bf}}{n_{bf}^{max}} \cdot \frac{1}{1 + I_{syn}^{int} + I_{syn}^{ord} + I_{sub}} \quad (5.7)$$

Poiché gli indici variazionali compaiono al denominatore, i valori di IP aumentano al diminuire dei valori dei tre indici (in questo modo, alti valori per IP implicano alta fissità). Allo stesso tempo le informazioni relative alle variazioni empiriche sono pesate allo stesso modo grazie all'operazione di somma. In questo modo un'espressione che presenti valori alti per uno solo degli indici può avere un valore risultante per IP simile a quello di un'espressione che presenti valori medi distribuiti su tutti gli indici variazionali. Tale struttura, quindi, permette di tenere in considerazione la flessibilità di comportamento delle espressioni terminologiche, quand'anche esse non risultino prototipiche e quindi fisse su tutti i fronti. Infine l'indice IP considera un fattore correttivo, dato dal rapporto di normalizzazione tra il numero di occorrenze di una data espressione e quello dell'espressione più frequente dell'insieme considerato (n_{bf}^{max}). Tale fattore, che risulta limitato tra 0 e 1, è necessario al fine di prendere in considerazione il fatto che basse occorrenze della forma base dell'espressione riducono l'affidabilità dei test empirici, in quanto la presenza o assenza di modificazioni non ha la possibilità di essere testata su un largo insieme di espressioni.

5.5 Analisi e risultati

Al fine di avere un'evidenza empirica delle performance dell'indice IP, si considera una lista di espressioni estratte dal corpus di fisica e composta da tutti i bigrammi relativi al pattern NA estratti automaticamente grazie al *tool mwetoolkit* (Ramisch *et al.*, 2010b). L'insieme di espressioni così composto annovera circa 22.700 entità.

Per ognuna di queste lo strumento computazionale svolge, quindi, i test di modifica variazionale e grazie ai valori ottenuti è possibile assegnare ad ogni espressione un valore di IP, che in questo modo si configura come una nuova misura d'associazione non basata su ipotesi statistiche, bensì puramente linguistiche.

Ordinando le espressioni secondo i valori di IP si ottiene una lista che appare analoga a quelle generalmente prodotte dalle misure d'associazione statistica, poiché l'indice risulta in grado di filtrare i candidati che non presentano un comportamento prototipico in riferimento alle caratteristiche sopra esposte e che per questo vengono spinte in fondo alla lista. Allo stesso tempo, le espressioni che ottengono valori molto alti per l'indice sembrano ben rappresentare i candidati più idonei ad essere classificati come polirematiche terminologiche. Le Tabelle 5.5 e 5.6 mostrano, rispettivamente, i primi e gli ultimi candidati della lista ordinata secondo i valori IP.

Rango	Espressione	IP
1	campo magnetico	0,9565
2	campo elettrico	0,6133
3	momento angolare	0,5717
4	meccanica quantistica	0,5205
5	calorimetro elettromagnetico	0,4748
6	modello standard	0,4259
7	valore medio	0,4206
8	massa invariante	0,3683
9	energia cinetica	0,3630
10	campo gravitazionale	0,3423
11	campo elettromagnetico	0,3314
12	relatività generale	0,3155
13	buco nero	0,2997
14	meccanica classica	0,2591
15	carica elettrica	0,2395

Tabella 5.5: Prime 15 espressioni della lista di bigrammi NA del corpus di fisica ordinate secondo l'indice IP.

Al fine di valutare le performance di estrazione dell'indice IP, si è scelto di comparare i risultati ottenuti sulla lista di candidati estratti dal corpus con le due ben

Rango	Espressione	IP
22.686	entità indipendente	0,000007
22.687	caso tale	0,000007
22.688	fotone due	0,000005
22.689	condizione fondamentale	0,000005
22.691	parte maggiore	0,000004
22.692	ambito magnetico	0,000004
22.693	condizione finale	0,000003
22.694	dimensione media	0,000003
22.695	forma standard	0,000002

Tabella 5.6: Ultime espressioni della lista di bigrammi NA del corpus di fisica ordinate secondo l'indice IP.

note misure di associazione statistica log-likelihood (Dunning, 1993) e Pointwise Mutual Information (Church & Hanks, 1990) che, per quanto detto nel paragrafo 2.4.3, quantificano due diversi aspetti dell'associazione fra parole e che quindi permettono di avere a disposizione due diverse prospettive delle strategie di estrazione.

Grazie al *tool* “mwetoolkit” è possibile assegnare ad ogni bigramma della lista un valore di Log-likelihood (LL) e Pointwise Mutual Information (PMI) e ordinare di conseguenza la lista di espressioni secondo l'uno o l'altro punteggio di associazione. A questo punto le performance di estrazione dell'indice IP e delle due misure di associazione vengono valutate sulla base del tasso di presenza di *true positives* man mano che la lista ordinata secondo ciascuno dei tre punteggi viene scorsa fino alla fine⁵. Sono considerate *true positives* tutte le espressioni estratte dal corpus che sono presenti nella lista di polirematiche fisiche del GRADIT, che viene quindi assunta a *gold standard* dell'esperimento.

La Figura 5.1 mostra le curve che rappresentano i tassi di estrazione di *true positives* per le tre misure.

Come si vede, LL e IP hanno performance simili rispetto all'identificazione di polirematiche, il che indica che i test semantici e sintattici su dati empirici sono in grado di fornire buoni risultati anche nell'individuazione automatica. Il risultato meno soddisfacente di PMI può essere giustificato dal fatto che non è stato applicato alcun filtro in frequenza all'insieme iniziale di espressioni (tutte le espressioni, cioè, sono state considerate indipendentemente dal loro numero di occorrenze) e tale misura d'associazione è nota per sovrastimare espressioni di bassa frequenza che sono spesso *false positives* (Evert, 2008).

⁵Considerando, nella valutazione, l'intero insieme di binomi NA del corpus, è possibile interpretare l'andamento della curva come la *precisione* delle misure al variare della copertura dei candidati.

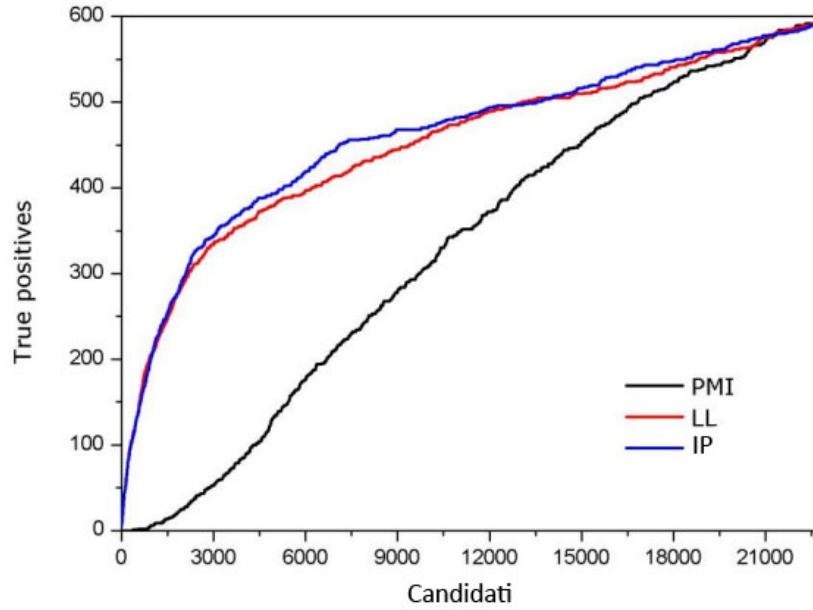


Figura 5.1: Confronto tra i tassi di estrazione di *true positives* di Pointwise Mutual Information (PMI, linea nera), Log-likelihood (LL, linea rossa) e Indice di Prototipicità (IP, linea blu).

Si può notare come per i primi 1.800 candidati (in corrispondenza di circa il 40% dei *true positives* estratti) LL ottenga risultati leggermente migliori rispetto all'indice IP. Tuttavia, per i restanti 20.900 candidati IP risulta quasi sempre il miglior estrattore. Ciò sembra indicare che su vasta scala l'indice IP sia maggiormente utile ai fini lessicografici o in studi dove risulti fondamentale l'estrazione del più alto numero possibile di polirematiche e non ci sia interesse solo verso il nucleo ristretto di espressioni che appaiono in cima alla lista.

A completamento dell'analisi, si è infine considerato anche un filtro in frequenza sulla lista dei bigrammi candidati, al fine di minimizzare i problemi legati ad espressioni di bassa frequenza, che interessano in particolare la misura PMI. Si genera, quindi, una nuova lista di espressioni che contiene solo bigrammi con frequenza di occorrenza $f \geq 30$ (per un totale di 301 espressioni), svolgendo la stessa procedura vista sopra.

Poiché un numero di occorrenze che superi 30 può essere considerato relativamente cospicuo e quindi affidabile in relazione ai test empirici svolti dallo strumento computazionale, si è scelto di considerare anche una variante "pura" dell'indice IP, che non viene corretta dall'informazione legata alla frequenza delle espressioni e che è calcolata secondo la seguente formula:

$$IP_p = \frac{1}{1 + I_{syn}^{int} + I_{syn}^{ord} + I_{sub}} \quad (5.8)$$

La Figura 5.2 mostra i tassi di individuazione di *true positives* per le quattro misure considerate.

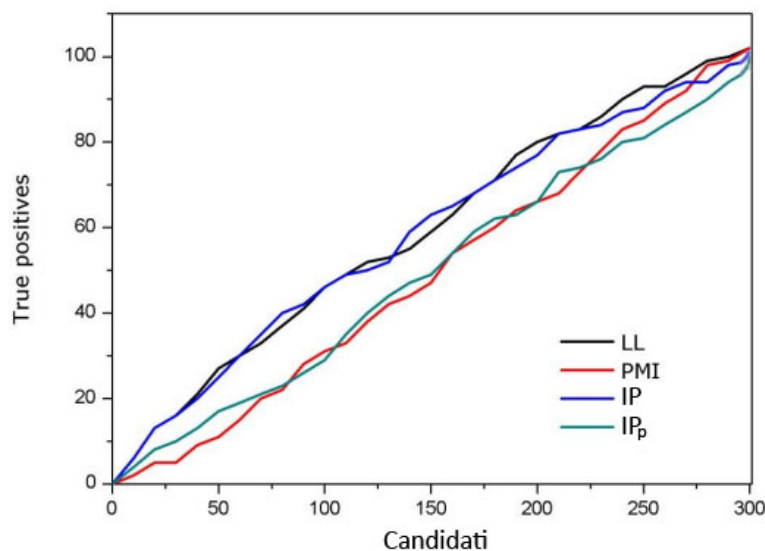


Figura 5.2: Confronto tra i tassi di estrazione di *true positives* di Log-likelihood (LL, linea nera), Pointwise Mutual Information (PMI, linea rossa), Indice di Prototipicità (IP, linea blu) e della versione pura dell'Indice di Prototipicità (IP_p, linea verde) su una lista di candidati con frequenza di occorrenza maggiore di 30.

Ancora una volta LL e IP si dimostrano la scelta migliore in termini di individuazione automatica e le loro performance risultano quasi identiche. PMI risulta, ora, maggiormente efficace, come c'era da aspettarsi, nonostante il suo tasso di estrazione rimanga meno efficiente di LL e IP. Infine la versione pura dell'indice di prototipicità IP_p mostra un tasso di estrazione che è chiaramente migliore rispetto a quello di PMI per i primi 80 candidati, mentre risulta comparabile ad esso per i restanti.

5.6 Conclusioni

Questo ristretto caso di studio sulla terminologia del settore fisico ha dimostrato come le caratteristiche variazionali sintagmatiche e paradigmatiche possono giocare un ruolo interessante nello studio e l'individuazione delle espressioni terminologiche, almeno in relazione al pattern NA.

I risultati indicano che anche nel caso della fisica, la terminologia tende a concentrarsi verso il polo delle polirematiche, mostrando quindi inibizione di variazioni sia sintagmatiche che paradigmatiche.

La comparabilità delle performance dell'indice IP con le misure d'associazione dimostra l'efficacia dell'attestazione delle variazioni in qualità di ulteriore metodologia per l'estrazione automatica, utile ad integrare i metodi statistici già disponibili.

Le ottime performance dell'indice IP su vasta scala possono essere motivate dal fatto che le proprietà sintattiche e semantiche, a differenza di quelle statistiche,

mostrano maggiore robustezza e affidabilità anche in presenza di espressioni meno frequenti.

L'evidenza che l'indice di prototipicità abbia performance migliori su vasta scala, inoltre, mostra come tale indice abbia potenzialità di impiego in particolare in lessicografia, dove si è in generale interessati al collezionare nel modo più efficiente il maggior numero di espressioni possibile, focalizzandosi sull'intero elenco di candidati e non solo sulle espressioni che appaiono in cima alla lista con il miglior punteggio.

Tuttavia, le performance ridotte della versione semplificata dell'indice IP, che non coinvolge l'informazione di frequenza, mostrano quanto il numero di occorrenze giochi inevitabilmente un ruolo di supporto e garantisca maggiore affidabilità all'individuazione della terminologia.

Conclusioni e lavori futuri

In questo lavoro si è provato a fornire una nuova prospettiva alla categorizzazione del vasto insieme delle espressioni multiparola italiane, spesso ricondotto alla dualità tra polirematiche e collocazioni. La carenza di studi che indagassero in maniera esplorativa ma sistematica corpora dell'italiano per analizzare il comportamento empirico di questo tipo di espressioni ha motivato la costruzione di uno strumento computazionale che fosse in grado di fornire statistiche sulle possibili variazioni sintagmatiche e paradigmatiche di combinazioni relative a pattern specificati. A valle della produzione dei dati quantitativi si è optato per un'analisi qualitativa delle espressioni che mirasse ad evidenziare caratteristiche omogenee a seconda della distribuzione di queste in relazione a presenza o assenza delle differenti variazioni e permettesse, in questo modo, di proporre la definizione di poli categoriali lungo un *continuum*.

Si è riusciti, in questo modo, a evidenziare come in generale la possibile categorizzazione cambi in dipendenza del pattern e della natura del sintagma analizzato.

Mentre le polirematiche, che individuano entità unitarie semanticamente o terminologiche, sono contraddistinte dall'impossibilità di modificazione indipendentemente dalla natura del sintagma, le collocazioni e gli abbinamenti preferenziali sembrano instaurarsi sotto la spinta di due possibili meccanismi opposti e complementari, contraddistinti dall'inibizione di una sola delle modificazioni sintattico-sintagmatiche e paradigmatiche. I sintagmi nominali preferiscono la prima opzione, mentre sintagmi in coordinazione e verbali (limitatamente, per ora, al VDN) la seconda.

La tendenziale fissità della terminologia polirematica è stata verificata, poi, su un caso di studio sul linguaggio tecnico-specialistico della fisica. Grazie ad un corpus costruito *ad hoc* e alle possibilità dello strumento computazionale, si sono analizzate le caratteristiche variazionali delle espressioni relative al pattern NA, grazie alle quali è stato possibile costruire un indice di prototipicità efficientemente utilizzato per l'individuazione automatica della terminologia. Le performance dell'indice sono state, infatti, confrontate con quelle di note misure d'associazione statistica, riscontrando comportamenti simili o migliori da parte della misura variazionale.

Tale valutazione è stata possibile grazie alla disponibilità di un *gold standard* utilizzabile poiché rappresentativo di un insieme di espressioni con uno statuto categoriale definito, a differenza di quanto accaduto per il linguaggio generale, per cui non si è riscontrata la presenza di risorse che garantissero un'affidabilità di discernimento categoriale tra le diverse tipologie di fenomeni multiparola.

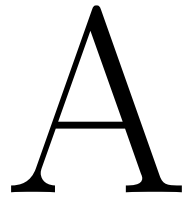
In tal senso, una delle principali linee future di indagine che il presente lavoro apre è un approfondimento empirico e lessicografico sul comportamento variazionale delle espressioni del linguaggio generale presenti nelle risorse disponibili per

l'italiano (De Mauro 1999-2007, Urzì 2009, Lo Cascio 2011, Tiberii 2012), in primis attraverso l'individuazione degli insiemi non intersecantisi tra polirematiche (in GRADIT) e collocazioni (nei dizionari combinatori) e lo studio del comportamento di tali espressioni anche in eventuali altri corpora di italiano.

Lo sviluppo dei test empirici per ulteriori pattern, specialmente verbali o a più componenti, risulta inoltre un passo obbligato nelle future versioni dello strumento computazionale. Oltre agli indirizzi suddetti, infatti, quest'ultimo ha la possibilità di essere migliorato su diversi fronti:

- attraverso l'inclusione o il confronto con ulteriori risorse da utilizzare come tesauri;
- attraverso la creazione di specifiche classi di entità omogenee quali specificatori temporali (es. nomi dei mesi), unità di misura, valute, ecc. che includono termini generalmente privi di sinonimi, ma sostituibili all'interno della stessa classe. I lemmi presenti in queste classi potrebbero costituire un ulteriore test di sostituzione da affiancare a quello sinonimico;
- attraverso la considerazione di test incrociati, ad esempio di sostituzione e interrompibilità allo stesso tempo, o di altre possibili combinazioni;
- attraverso il riconoscimento automatico di espressioni comprese entro locuzioni più ampie (queste ultime individuate attraverso criteri statistici o di semplice frequenza).

Le conclusioni sopra delineate, nonché i propositi di sviluppo della metodologia di ricerca e degli strumenti mostrano quanto ancora resti da indagare ai fini di una sufficiente comprensione dei fenomeni multiparola dell'italiano. Nondimeno, il presente lavoro è riuscito a mostrare la validità di una linea di approfondimento empirico che merita attenzione nel vasto panorama di studi sull'argomento.



Part of speech Tagset - PAISÀ

Si riporta nel seguito, così come fornita dagli sviluppatori, la documentazione relativa alla descrizione dell'insieme di etichette grammaticali generali (*coarse*) e specifiche (*fine-grained*) utilizzate per il corpus PAISÀ, disponibile anche all'indirizzo

http://www.corpusitaliano.it/static/documents/POS_ISST-TANL-tagset-web.pdf

ISST-TANL Tagsets

The morpho-syntactic and dependency tagsets were jointly developed by the Istituto di Linguistica Computazionale “Antonio Zampolli” (ILC-CNR) and the University of Pisa in the framework of the TANL (Text Analytics and Natural Language processing) project and was used for the annotation of the ISST-TANL dependency annotated corpus.

ISST-TANL morpho-syntactic tagset

The ISST-TANL part-of-speech tags are based on the ILC/PAROLE tagset and are conformant to the EAGLES international standard. The table below documents the 14 coarse-grained pos tags (column 1) and the 37 fine-grained tags (column 2) used for ISST-TANL annotation.

Coarse-grained tag	Fine-grained tag	Description	Examples	Contexts of use
A	A	adjective	<i>bello, buono, pauroso, ottimo</i>	una bella passeggiata un ottimo attaccante una persona paurosa
A	AP	possessive adjective	<i>mio, tuo, nostro, loro</i>	a mio parere il tuo libro
B	B	adverb	<i>bene, fortemente, malissimo, domani</i>	arrivo domani sto bene
B	BN	negation adverb	<i>non</i>	non sto bene
C	CC	coordinative conjunction	<i>e, o, ma, ovvero</i>	i libri e i quaderni vengo ma non rimango
C	CS	subordinative conjunction	<i>mentre, quando</i>	quando ho finito vengo mentre scrivevo ho finito l'inchiostro
D	DE	exclamative determiner	<i>che, quale, quanto</i>	che disastro! quale catastrofe!
D	DI	indefinite determiner	<i>alcuno, certo, tale, parecchio, qualsiasi</i>	alcune telefonate parecchi giornali qualsiasi persona
D	DQ	interrogative determiner	<i>che, quale, quanto</i>	che cosa quanta strada quale formazione
D	DR	relative determiner	<i>cui, quale</i>	i cui libri
D	DD	demonstrative	<i>questo, codesto,</i>	questo denaro

		determiner	<i>quello</i>	quella famiglia
E	E	preposition	<i>di, a, da, in, su, attraverso, verso, prima_di</i>	a casa del poeta prima_di giorno verso sera
E	EA	articulated preposition	<i>del, alla, dei, nelle</i>	nella casa il prezzo del pane
F	FB	“balanced” punctuation	() “ ” ‘ ’ - -	il gatto – che conoscete –
F	FC	clause boundary punctuation	, ;	ha detto : Vieni!
F	FF	comma, hyphen	,	mele, pere e banane due-trecento persone
F	FS	sentence boundary punctuation	. ? !	mele, pere e banane. cosa vuoi?
I	I	interjection	<i>ahimè, beh, ecco, grazie</i>	Beh , che vuoi?
N	N	cardinal number	<i>uno, due, cento, mille, 28, 2000</i>	due partite 28 anni
N	NO	ordinal number	<i>primo, secondo, centesimo</i>	secondo posto
P	PD	demonstrative pronoun	<i>questo, quello, costui</i>	quello di Roma costui uccide
P	PE	personal pronoun	<i>egli, lui, esso noialtri, voialtri, essi io, me, tu, te</i>	io parto lo mangio
P	PI	indefinite pronoun	<i>chiunque, ognuno, molto</i>	chiunque venga i diritti di ognuno
P	PP	possessive pronoun	<i>mio, tuo, suo, loro, proprio</i>	il mio è qui più bella della loro
P	PQ	interrogative pronoun	<i>che, chi, quanto</i>	non so chi parta quanto costa? che ha fatto ieri?
P	PR	relative pronoun	<i>che, cui, quale</i>	ciò che dice il quale afferma a cui parlo
P	PC	clitic pronoun	<i>ci, vi, mi, ti, la, le</i>	lo vidi li ho sentiti aver la

				le dissero, le videro mi dicono ci sposiamo vi credo si sente, si sentono ci vado spesso
R	RD	determinative article	<i>il, lo, la, i, gli, le</i>	il libro i gatti
R	RI	indeterminative article	<i>uno, un, una</i>	un amico una bambina
S	S	common noun	<i>amico, insegnante, verità</i>	l' amico la verità
S	SA	abbreviation	<i>ndr, a.C., d.o.c., km</i>	30 km sesto secolo a.C.
S	SP	proper noun	<i>Monica, Pisa, Fiat, Sardegna</i>	Monica scrive
T	T	predeterminer	<i>tutto, entrambi, ambedue</i>	tutte le notizie ambedue le idee
V	VA	auxiliary verb	<i>avere, essere, venire</i>	il peggio è passato ho scritto una lettera viene fatto domani
V	VM	modal verb	<i>volere, potere, dovere, solere</i>	non posso venire vuole il libro
V	V	main verb	<i>mangio, avere, passato, camminando</i>	il peggio è passato ho scritto una lettera vengo domani
X	X	residual class	it includes formulae, unclassified words, alphabetic symbols and the like	distanziare di 43'' mi piacce

B

Dati sul pattern NA

Si riporta di seguito la lista, in ordine alfabetico, delle 500 espressioni relative al pattern NA estratte in PAISÀ e i dati relativi al comportamento variazionale di ognuna delle espressioni.

acqua caldo 840	atletica leggero 1760
acqua dolce 1789	attività agonistico 969
acqua potabile 819	attività agricolo 782
adattamento cinematografico 785	attività commerciale 1178
aereo militare 2808	attività economico 1206
aeronautica militare 754	attività politico 1450
altare maggiore 3117	attività produttivo 787
ambiente naturale 769	attività sportivo 952
amministratore delegato 1082	attore teatrale 756
amministrazione comunale 1515	autorità competente 780
anidride carbonico 1382	azienda agricolo 777
anno consecutivo 997	azienda italiano 906
anno precedente 4271	a?? a?? 1042
anno scorso 2956	banda largo 1123
anno seguente 7296	base aereo 803
anno successivo 18341	base militare 937
architettura civile 737	bene culturale 735
architettura religioso 1127	braccio destro 1298
area geografico 771	brano musicale 1495
area metropolitano 1473	buco nero 1583
area naturale 1322	calcio italiano 1298
area protetto 1085	cambiamento climatico 915
area urbano 1145	cambio automatico 906
argomento simili 3397	campagna elettorale 2691
arma nucleare 982	campagna militare 1301
arte contemporaneo 1120	campagna pubblicitario 1353
arte marziale 3369	campionato europeo 2407
articolo simile 4800	campionato italiano 3404

- campionato mondiale 3760
campionato nazionale 2025
campo elettrico 866
campo magnetico 2417
canale televisivo 948
cantiere navale 1003
cappella laterale 1002
cappucci italiano 792
caratteristica fisico 2825
caratteristica peculiare 744
caratteristica principale 1312
caratteristica tecnico 2514
carriera militare 1124
carriera musicale 883
carriera politico 1367
carriera solista 1054
carro armato 3224
cartone animato 3031
casa automobilistico 1230
casa discografico 4041
casa editore 5148
caso contrario 843
caso particolare 1327
catena montuoso 1571
cellula staminali 1474
cemento armato 956
centrale nucleare 1015
centro cittadino 1421
centro commerciale 2211
centro culturale 782
centro sociale 1991
centro storico 6787
centro urbano 1580
chiesa cattolico 7127
chiesa parrocchiale 3221
chilometro quadrato 805
chitarra elettrico 948
cielo buio 967
cielo limpido 764
cielo serale 2603
cinema italiano 1264
cinta murario 1838
circolo polare 1583
città italiano 1380
cittadino italiano 812
classe dirigente 946
classe operaio 958
classe politico 847
classe sociale 1134
classe spettrale 3100
classifica finale 1572
classifica generale 1132
codice civile 825
codice penale 919
collegamento esterno 131642
colonna sonoro 9197
colore bianco 1106
colore giallo 780
colore nero 817
colore rosso 1323
colore verde 906
commento anonimo 4382
compagnia aereo 2410
compagnia teatrale 753
competizione internazionale 1680
competizione nazionale 1555
compositore classico 774
comunità ebraico 1165
comunità internazionale 736
comunità scientifico 998
concessionario esclusivo 1228
condizione ambientale 774
condizione climatico 797
condizione economico 908
conflitto mondiale 1874
consigliere comunale 1346
consigliere regionale 767
consiglio comunale 1628
consorzio agrario 1238
contenuto video 2796
contestato storico 834
coppia massimo 1161
corpo celeste 803
corpo umano 1312
costa occidentale 824
costa orientale 884

-
- crescita economico 853
crisi economico 2075
cultura popolare 1521
dato personale 1089
decennio successivo 1256
decreto legislativo 887
difesa aereo 861
difficoltà economico 1122
direttore artistico 1218
direttore generale 2288
diritto canonico 740
diritto civile 2281
diritto internazionale 1027
diritto umano 3081
drama politico 944
edificio religioso 1680
edificio sacro 886
edizione italiano 1476
effetto collaterale 1790
effetto speciale 1701
elezione europeo 846
elezione politico 2647
elezione presidenziale 1540
elezione regionale 813
emisfero boreale 747
emisfero celeste 2397
emittente televisivo 969
energia cinetico 923
energia elettrico 2535
ente locale 1220
ente pubblico 1303
epoca medievale 971
epoca romano 2680
equatore celeste 2050
equazione differenziale 825
esercito francese 1112
esercito romano 938
esercito tedesco 794
essere umano 6704
essere vivente 1493
età adulto 942
età medio 1125
età moderno 821
etichetta discografico 1780
evoluzione demografico 5748
facezia filosofico 739
fama internazionale 897
famiglia nobile 1192
famiglia reale 1608
fascia tropicale 1070
fase finale 2349
fase iniziale 820
figlio illegittimo 1051
figlio maggiore 1290
film drammatico 2536
fonte primario 881
fonte storico 775
fonte video 2742
forza aereo 1108
forza armato 3882
forza militare 1040
forza politico 1365
fratello maggiore 2527
fratello minore 1120
galleria fotografico 3082
gas naturale 863
genere letterario 761
genere musicale 2048
genere umano 910
gente comune 847
geografia fisico 1652
gioco olimpico 754
giorno precedente 907
giorno scorso 1143
giorno seguente 1995
giorno successivo 3026
governo italiano 1568
greco antico 880
gruppo etnico 2215
gruppo montuoso 820
gruppo musicale 5025
guerra civile 6426
guerra freddo 1643
guerra mondiale 27803
impatto ambientale 1063
impegno politico 772

imperatore bizantino 796	modo migliore 791
imperatore romano 927	modo particolare 1306
impero romano 1069	modo tale 2224
impianto sportivo 999	mondo esterno 735
impiego operativo 737	mondo intero 1106
incidente stradale 1737	mondo occidentale 735
informazione pubblicitario 1219	mondo reale 889
intelligenza artificiale 942	movimento politico 1048
intervento chirurgico 1143	musica classico 1428
invasione barbarico 779	musica elettronico 1060
istituto religioso 1079	musica popolare 755
lato destro 1123	nativo americano 893
lato sinistro 1113	navata centrale 1243
legge elettorale 805	navata laterale 743
letteratura italiano 743	navata unico 747
linea ferroviario 3749	nome attuale 1025
lingua francese 1081	nome comune 932
lingua inglese 2998	nome originale 844
lingua italiano 3500	nucleo familiare 848
lingua latino 891	numero reale 1249
lingua originale 855	occhio nudo 2630
lingua tedesco 1528	opera letterario 1681
lingua ufficiale 1388	opera principale 1462
livello internazionale 2678	opera pubblico 791
livello mondiale 1880	opera teatrale 2264
livello nazionale 2618	operazione militare 1120
livello superiore 784	opinione pubblico 3916
luogo comune 1172	ordine cronologico 1044
macchina fotografico 1038	ordine pubblico 1423
magnitudine assoluto 2139	ordine religioso 835
magnitudine pari 2284	organizzazione internazionale 894
mano destro 1244	oro olimpico 879
mano sinistro 943	paese europeo 1903
marmo bianco 833	paese occidentale 766
mercato italiano 920	pallavolo femminile 834
mese estivo 736	pannello solare 763
mese successivo 1200	parco nazionale 1501
messaggio privato 4275	parlamento europeo 1488
metro quadrato 1213	parte alto 1244
metro quadro 817	parte anteriore 1274
mezzo pubblico 774	parte basso 928
miniserie televisivo 833	parte centrale 2169
modo diverso 1419	parte finale 940

-
- parte inferiore 2094
 - parte integrante 2586
 - parte meridionale 1731
 - parte occidentale 1511
 - parte orientale 1523
 - parte posteriore 1701
 - parte settentrionale 1649
 - parte superiore 2934
 - partito comunista 1328
 - partito politico 3142
 - penisola iberico 809
 - periodo estivo 1258
 - periodo migliore 2655
 - periodo storico 1617
 - periodo successivo 737
 - persona morto 862
 - personaggio famoso 750
 - personaggio immaginario 1027
 - personaggio principale 2091
 - personalità legato 17233
 - piano nobile 814
 - piano superiore 1582
 - pianta erbaceo 747
 - pianta rettangolare 749
 - piazza principale 768
 - pietra miliare 953
 - pinna dorsale 873
 - pista ciclabile 744
 - politica economico 948
 - politica estero 2414
 - politica interno 768
 - popolazione civile 740
 - popolazione locale 1701
 - posizione geografico 866
 - posizione strategico 1051
 - potenza massimo 2059
 - potere politico 1231
 - premio letterario 843
 - prima time 2664
 - problema economico 824
 - problema tecnico 747
 - processo produttivo 880
 - prodotto agricolo 744
 - prodotto tipico 947
 - produzione artistico 770
 - produzione cinematografico 746
 - professore ordinario 750
 - progetto speciale 1281
 - programma televisivo 3180
 - proprietà privato 1166
 - punto debole 1001
 - quartier generale 2047
 - rapporto sessuale 1177
 - regime fascista 1224
 - regione ecclesiastico 827
 - regione italiano 747
 - regione temperato 1176
 - reperto archeologico 923
 - rete ferroviario 758
 - rete televisivo 1388
 - ricerca scientifico 2002
 - riserva naturale 1064
 - risorsa naturale 890
 - risultato finale 1216
 - riva destro 744
 - ruolo fondamentale 1261
 - ruolo importante 2352
 - sala cinematografico 1118
 - scavo archeologico 980
 - scena musicale 768
 - sci alpino 1097
 - scienza naturale 742
 - scuola elementare 2044
 - scuola medio 1293
 - scuola pubblico 871
 - scuola superiore 1715
 - secolo scorso 2190
 - secolo successivo 2508
 - secolo x 1199
 - secondo singolo 1415
 - sede vescovile 1854
 - segretario generale 1387
 - senso stretto 849
 - sequenza principale 2888
 - serie animato 2048
 - serie televisivo 9064

servizio militare 1787	strada principale 793
servizio pubblico 1743	strada provinciale 1508
servizio segreto 2529	strada statale 2521
settimana scorso 807	strumento musicale 1723
settimana successivo 817	studente universitario 772
settore giovanile 998	successo commerciale 1834
sistema economico 776	successo internazionale 814
sistema nervoso 1649	sviluppo economico 1711
sistema operativo 6416	target comm.le 2014
sistema politico 797	telefono cellulare 1206
sistema solare 5982	televisione digitale 740
sistema stellare 1020	tema principale 741
sito archeologico 2816	temperatura medio 1558
sito ufficiale 2393	tempo indeterminato 741
situazione economico 816	tempo libero 1684
situazione politico 901	tempo necessario 871
società calcistico 2107	tempo pieno 1358
società civile 1327	tempo reale 1667
società italiano 1038	tempo recente 1135
software libero 1346	tempo stesso 2772
sorella maggiore 973	tempo supplementare 964
sorella minore 795	territorio comunale 4564
spazio vettoriale 833	territorio italiano 1221
spettacolo teatrale 1265	territorio nazionale 1613
spot pubblicitario 889	testamento biologico 1126
squadra nazionale 758	testata giornalistico 1001
stagione estivo 814	titolo mondiale 1547
stagione precedente 794	titolo nazionale 1044
stagione regolare 956	titolo originale 1947
stagione seguente 776	torre campanario 941
stagione successivo 3289	tradizione popolare 1114
stato attuale 1109	traduzione italiano 765
stato italiano 1622	trasmissione televisivo 2019
stato maggiore 1096	trasporto pubblico 2158
stazione ferroviario 3928	trazione integrale 878
stazione meteorologico 1305	truppa francese 855
stazione spaziale 991	truppa tedesco 835
stella gigante 1173	uomo politico 1102
stile barocco 1009	uso comune 1184
stile gotico 787	valor militare 1187
stile libero 740	valore assoluto 1264
stile musicale 882	velocità massimo 3270
stile romanico 777	velocità radiale 2220

versione italiano 1812	vita quotidiano 2278
versione originale 2042	vita reale 737
via libero 862	vita sociale 1317
vicariato apostolico 1237	vita umano 1123
video musicale 1411	vittoria finale 833
vita politico 1567	zona industriale 809
vita privato 3536	
vita pubblico 755	

Di seguito viene inoltre riportato l'insieme dei dati prodotti dallo strumento computazionale per tutte le espressioni considerate. I nomi delle colonne sono precisati dalla seguente legenda:

ESPRESSIONE	Espressione considerata
FREQ	Frequenza di occorrenza (componenti lemmatizzati)
FREQ_INT	Numero di occorrenze dell'espressione interrotta
INTERROMP	Valore di I_{syn}^{int}
FREQ_REV	Numero di occorrenze dell'espressione con componenti invertiti
ORDINE INVERSO	Valore di I_{syn}^{ord}
MOD_SINTAG	Valore di I_{syn}
FREQ_SIN1	Numero di occorrenze delle espressioni con primo componente sostituito (valore non intero poiché normalizzato in base alla similarità semantica)
DISP1	Dispersione della sostituzione del primo componente
FREQ_SIN2	Numero di occorrenze delle espressioni con secondo componente sostituito (valore non intero poiché normalizzato in base alla similarità semantica)
DISP2	Dispersione della sostituzione del secondo componente
MOD_PARAD	Valore di I_{sub}
FORMA_PREVAL	Forma con il più alto numero di occorrenze
FREQ_P	Numero di occorrenze della forma prevalente
FORMA_PREVAL_2	Forma con il secondo numero di più alto di occorrenze
FREQ_P_2	Numero di occorrenze della seconda forma prevalente
MOD_FLESS	Valore di I_{infl}
ERR	Se contrassegnato da X, espressione scartata per errori dovuti al corpus

emissario boreale	7477	913	0.55	0	0.55	0	0.511756545	0	0	0	1.4343751359	0.0019165009	emissario boreale	7477	0	1	0.0013386881	0
pianta erbacea	7477	2	0.002670227	781	0	0.5117647059	0	0	0	0	0.7171875019	0	pianta erbacea	7477	0	262	0.3574297189	0
pianta erbacea	7477	78	0.0945454545	0	0.0945454545	0	0.0945454545	0	0	0	0	0.0009591696	pianta erbacea	7477	0	162	0.2275799746	0
regione italiana	7477	78	0.0945454545	0	0.0945454545	1333	1840684096	0.0068647795	0	0	0.166698352	0.6412105068	regioni italiane	7477	0	165	0.2208835341	0
problema tecnico	7477	12	0.0158102767	0	0.0158102767	72	12462316764	0.2590314339	0	0	20.877264445	0.1103227834	problemi tecnici	7477	0	582	0.2208835341	0
produzione cinematografica	746	6	0.007987234	0	0.007987234	243	6544246014	0.490902271	0	0	3.9324134214	0.2491847662	produzione cinematografica	557	0	178	0.2533512064	0
caratteristica peculiare	744	89	0.0068427371	122	0.0068427371	73	0696790254	0.4850593231	0	0	683.3806237489	0.0337757162	caratteristiche peculiari	477	0	302	0.4059139785	0
prodotto agricolo	744	3	0.0040160643	1	0.0040160643	49	7259518547	0.4135101496	0	0	1.6600856639	0.0646051056	prodotti agricoli	697	0	45	0.063172043	0
pianta ciclabile	744	3	0.0040160643	0	0.0040160643	44	9387151772	0.0547319349	0	0	0.5569609709	0.0569609709	piante ciclabili	371	0	369	0.501344086	0
rivera destra	744	3	0.0040160643	0	0.0040160643	391	2483482514	0.0147766753	0	0	0.4621988391	0.3449035039	rivera destra	744	0	0	0.50269179	0
letteratura italiana	743	44	0.05509085133	0	0.05509085133	1.6609011137	0	0.0147766753	0	0	0.2022304127	0.0022304127	letteratura italiana	741	0	1	0.00269179	0
lavorata laterale	743	13	0.0171957672	0	0.0171957672	0	0	0	0	0	2.0761264172	0	lavorata laterale	611	0	129	0.1778581427	0
scienza naturale	742	4	0.0035013303	0	0.0035013303	6.0950455195	1.0389207862	0	0	0	12.5499701975	0.1205822735	scienza naturale	693	0	41	0.003423181	0
tema principale	742	49	0.062052165	64	0.062052165	300.4217791384	0.13596000432	105.2641004267	0	0	103.656810761	0.016990486	temi principali	520	0	221	0.296245614	0
tempo indeterminato	741	5	0.0067024129	0	0.0067024129	8.9982433102	0.2175664451	0	0	0	0.12869899077	0.12869899077	tempi indeterminati	73	0	5	0.0094466937	0
diritto canonico	740	11	0.0146471372	0	0.0146471372	28.7378707977	0.2141918237	0	0	0	0.1970849163	0.0633126457	diritto canonico	739	0	1	0.0013513514	0
stile libero	740	10	0.0133333333	0	0.0133333333	123.0831498218	0.0620695322	0	0	0	11.5296255288	0.9175559018	stili liberi	739	0	1	0.0013513514	0
popolazione civile	740	4	0.0053763441	0	0.0053763441	8.0557965703	0.5417381422	0	0	0	3.1869943826	0.0149655681	popolazioni civili	655	0	84	0.1148648849	0
televisione digitale	740	0	0	0	0	0	0	0.0569147149	0	0	0	0.0538539072	televisione digitale	732	0	8	0.0108108108	0
faccia filosofica	739	0	0	0	0	0	0	0	0	0	0	0	faccia filosofica	739	0	0	0	0
tempo successivo	737	194	0.2043780881	112	0.1319199058	0.29339844679	31.878205019669	0.0035848319	0	0	0.0953129487	0.8160482542	periodi successivi	6	0	136	0.1858887381	0
vita reale	737	17	0.0256464191	8	0.0107388255	0.0328008398	172.5500992115	0.1453566492	24.3437294589	0	0.75028029852	0.2108358288	vita reale	733	0	2	0.0054274084	0
architettura civile	737	4	0.005981107	0	0.005981107	260.9363675812	0.0419306981	39.1126625539	0.0869788724	0	0.2892610044	0.2892610044	architetture civili	569	0	168	0.2279511533	0
impiego operativo	737	2	0.007063599	0	0.007063599	25.7995135189	0.8023324733	28.3213541857	0.0313365604	0	0.0660019395	0.0660019395	impiego operativo	737	0	0	0	0
comunità internazionale	736	120	0.1401869159	0	0.1401869159	1123.261641258	0.19124654231	37.5033032706	0.1924654231	0	0.6119780797	0.6119780797	comunità internazionali	730	0	6	0.0095108096	0
messe estivo	736	2	0.0027100271	0	0.0027100271	0	0	0	0.031175043	0	0.2205735748	0.2205735748	72mensa estivo	729	0	1	0.0013605442	0
mondo occidentale	735	27	0.0354330709	0	0.0354330709	208.001130364	0.0420170285	0	0	0	0.0680065644	0.0680065644	mondo occidentale	734	0	7	0.0095108096	0
bene culturale	735	2	0.016064257	0	0.016064257	233.2849838257	0.0339974213	0	0	0	0.2422169531	0.2422169531	beni culturali	699	0	35	0.4889795918	0
mondo esterno	735	3	0.0040650407	0	0.0040650407	445.3147734938	0.0325325056	51.129690098	0.3502792991	0	0.4031399489	0.4031399489	mondo esterno	730	0	5	0.0068072711	0

C

Dati sul pattern AN

Si riporta di seguito la lista, in ordine alfabetico, delle 500 espressioni relative al pattern AN estratte in PAISÀ e i dati relativi al comportamento variazionale di ognuna delle espressioni.

- commento 525	altro donna 732
alterno vicenda 526	altro elemento 703
alto carica 566	altro esempio 714
alto definizione 898	altro figlio 742
alto frequenza 537	altro film 667
alto grado 599	altro forma 684
alto livello 3671	altro genere 509
alto montagna 508	altro giocatore 550
alto numero 619	altro giorno 628
alto pressione 649	altro gruppo 1081
alto prestazione 686	altro lingua 651
alto qualità 1144	altro materiale 515
alto quota 1347	altro membro 1656
alto risoluzione 507	altro modo 673
alto tasso 565	altro mondo 528
alto temperatura 998	altro nazione 561
alto tensione 543	altro nome 517
alto valle 647	altro opera 900
alto valore 556	altro paese 1598
alto velocità 2806	altro parte 3789
alto voce 523	altro persona 1213
altro artista 500	altro personaggio 1188
altro aspetto 625	altro pianeta 497
altro caratteristica 568	altro punto 575
altro caso 690	altro ragazzo 736
altro città 900	altro sistema 543
altro componenti 556	altro specie 960
altro cosa 3062	altro squadra 727

altro storia 557	buono fede 941
altro tipo 1452	buono notizia 544
altro uomo 773	buono parte 6488
altro utente 4507	buono prestazione 545
altro versione 999	buono qualità 635
altro volta 868	buono rapporto 693
ampio spazio 823	buono risultato 1151
antico chiesa 1163	buono volontà 545
antico città 1373	caro amico 870
antico famiglia 562	certo importanza 620
antico nome 550	certo numero 3028
antico tradizione 913	certo punto 3183
antico via 499	certo tipo 573
attuale chiesa 561	comune italiano 2259
attuale presidente 506	discreto successo 1547
attuale via 662	diverso tempo 736
basso costo 1569	diverso tipo 783
basso livello 807	drammatico - 545
basso quota 1224	duro colpo 570
basso temperatura 604	duro prova 588
basso velocità 537	elevato numero 688
bello arte 509	enorme quantità 603
bello donna 775	enorme successo 1438
bello mezzo 728	estremo destra 1385
bello pò 592	estremo sinistra 674
bello ragazzo 875	eventuale richiesta 2774
bravo ragazzo 531	ex compagno 862
breve distanza 1111	ex marito 628
breve durata 974	ex membro 766
breve parentesi 574	ex ministro 667
breve periodo 4279	ex moglie 570
breve storia 511	ex presidente 1076
breve tempo 4345	facente parte 2439
breve termine 720	fine stagione 697
breve tratto 498	fondamentale importanza 648
buon fine 719	forte influenza 592
buon livello 492	forte legame 593
buon numero 912	giovane artista 585
buon occhio 539	giovane donna 1355
buon risultato 701	giovane età 1338
buon senso 989	giovane ragazzo 890
buon successo 1860	gran maestro 268
buono domenica 684	gran numero 4515

-
- gran parte 20001
gran premio 1041
gran quantità 574
gran voce 600
grande amico 1208
grande amore 712
grande artista 785
grande attenzione 604
grande attore 552
grande azienda 738
grande battaglia 498
grande capacità 788
grande centro 682
grande città 2091
grande difficoltà 503
grande dimensione 1973
grande distanza 645
grande evento 717
grande famiglia 667
grande forza 586
grande fratello 1339
grande gruppo 782
grande guerra 757
grande importanza 2192
grande impresa 590
grande influenza 615
grande interesse 1293
grande maestro 912
grande magazzino 509
grande maggioranza 687
grande manifestazione 506
grande nome 527
grande numero 1221
grande opera 1085
grande passione 811
grande popolarità 516
grande potenza 1028
grande pubblico 1069
grande quantità 2719
grande schermo 1713
grande successo 8694
grande talento 583
grande valore 1134
grande varietà 1157
grave crisi 829
grave danno 1187
grave incidente 699
grave infortunio 517
grave malattia 660
grave perdita 569
grave problema 1008
grosso modo 591
grosso problema 637
immediato dopoguerra 714
immediato vicinanza 733
importante centro 1727
importante città 583
importante contributo 500
importante opera 677
importante ruolo 522
inferiore rispetto 577
intenso attività 1394
intero area 581
intero città 578
intero famiglia 575
intero popolazione 562
intero regione 524
intero sistema 768
intero territorio 821
italiano - 1950
largo parte 1215
largo scala 1110
lato opposto 665
libero arbitrio 505
libero mercato 600
lieto fine 585
lungo anno 608
lungo carriera 734
lungo distanza 822
lungo durata 684
lungo periodo 3671
lungo raggio 1101
lungo serie 1511
lungo storia 587
lungo tempo 2897
lungo termine 1576

lungo viaggio 932
maggior numero 3237
maggior parte 31104
maggior ragione 566
maggior successo 2142
maggior attenzione 504
maggior città 528
maggior dimensione 559
maggior esponenti 500
maggior età 539
maggior importanza 504
maggior informazione 536
maggior rispetto 663
maggior successo 1003
massimo campionato 733
massimo divisione 669
massimo livello 805
massimo parte 839
massimo serie 3003
massimo splendore 532
medio dimensione 526
mezz' ora 1723
mezzo secolo 549
miglior attore 1959
miglior film 1586
miglior giocatore 649
miglior regista 351
miglior risultato 927
migliore amico 2171
migliore attore 2118
migliore giocatore 626
migliore risultato 517
minimo parte 606
nobile famiglia 1650
notevole importanza 1021
notevole interesse 770
notevole successo 1323
notevole sviluppo 491
numeroso opera 986
numeroso premio 561
numeroso specie 972
nuovo album 3391
nuovo allenatore 507
nuovo arma 551
nuovo avventura 577
nuovo canzone 651
nuovo casa 729
nuovo centro 537
nuovo chiesa 1248
nuovo città 667
nuovo classe 578
nuovo concezione 533
nuovo corso 586
nuovo costituzione 550
nuovo costruzione 597
nuovo disco 829
nuovo edificio 862
nuovo edizione 1122
nuovo fase 548
nuovo film 747
nuovo forma 1367
nuovo formazione 842
nuovo generazione 2720
nuovo governo 1570
nuovo gruppo 1078
nuovo idea 549
nuovo imperatore 501
nuovo impianto 568
nuovo lavoro 792
nuovo legge 974
nuovo linea 1065
nuovo materiale 629
nuovo millennio 692
nuovo modello 2046
nuovo modo 737
nuovo mondo 630
nuovo motore 1349
nuovo nome 863
nuovo ordine 633
nuovo parola 990
nuovo partito 658
nuovo personaggio 1366
nuovo presidente 692
nuovo prodotto 509
nuovo progetto 1469
nuovo programma 770

-
- nuovo provincia 660
 - nuovo quartiere 523
 - nuovo re 839
 - nuovo realtà 697
 - nuovo record 782
 - nuovo regola 501
 - nuovo sede 1055
 - nuovo serie 1777
 - nuovo sistema 2052
 - nuovo società 1009
 - nuovo squadra 645
 - nuovo stadio 540
 - nuovo stagione 907
 - nuovo stato 893
 - nuovo stazione 575
 - nuovo stile 607
 - nuovo strada 730
 - nuovo struttura 670
 - nuovo tecnica 758
 - nuovo tecnologia 1856
 - nuovo telefilm 1036
 - nuovo tipo 1240
 - nuovo versione 2761
 - nuovo vita 1143
 - omonimo film 640
 - omonimo romanzo 809
 - ottavo posto 596
 - ottimo ascolto 528
 - ottimo prestazione 523
 - ottimo risultato 1822
 - pari merito 607
 - pari opportunità 552
 - pari passo 686
 - particolar modo 631
 - particolare attenzione 1448
 - particolare importanza 827
 - particolare interesse 977
 - particolare riferimento 489
 - particolare rilievo 497
 - personal computer 913
 - piccolo borgo 717
 - piccolo centro 1121
 - piccolo chiesa 852
 - piccolo città 739
 - piccolo comunità 703
 - piccolo dimensione 2928
 - piccolo gruppo 1588
 - piccolo isola 671
 - piccolo numero 541
 - piccolo paese 963
 - piccolo parte 1878
 - piccolo quantità 841
 - piccolo ruolo 596
 - piccolo schermo 952
 - piccolo villaggio 1050
 - pieno titolo 782
 - poco distanza 1056
 - potente famiglia 575
 - presente visualizzazione 2742
 - principale attività 492
 - principale caratteristica 521
 - principale centro 849
 - principale città 1061
 - principale fonte 727
 - pronto soccorso 1034
 - prossimo anno 1410
 - prossimo giorno 734
 - prossimo mese 615
 - prossimo settimana 863
 - prossimo stagione 943
 - prossimo volta 7058
 - pubblico amministrazione 1416
 - pubblico dominio 558
 - pubblico istruzione 705
 - pubblico ministero 864
 - pubblico sicurezza 489
 - raro caso 694
 - recente studio 760
 - remoto server 924
 - ricco famiglia 507
 - rispettivo autore 494
 - santo - 706
 - scarso risultato 549
 - scarso successo 922
 - scorso anno 2336
 - scorso settimana 1061

serio problema 623
singolo estratto 510
singolo individuo 585
social network 1124
solo anno 595
solo punto 567
solo scopo 1013
solo stagione 569
solo volta 889
special modo 796
stessa città 693
stessa cosa 2030
stessa famiglia 627
stessa persona 556
stessa sorte 514
stessi anno 1088
stesso anno 23160
stesso autore 845
stesso discorso 298
stesso giorno 2081
stesso livello 585
stesso mese 911
stesso modo 4470
stesso motivo 554
stesso nome 1670
stesso numero 607
stesso periodo 4671
stesso piano 552
stesso sesso 594
stesso tempo 5668
stesso tipo 785
stesso titolo 573
stragrande maggioranza 1230
stretto contatto 805
stretto legame 524
stretto rapporto 570
tardo estate 554
tardo età 493
tenero età 869
time longer 678
tipico esempio 519
tutt' oggi 1924
tutt' ora 1984

ulteriore informazione 544
ultimo album 1336
ultimo anno 14021
ultimo apparizione 687
ultimo capitolo 511
ultimo caso 951
ultimo decennio 2413
ultimo edizione 917
ultimo episodio 1360
ultimo fase 895
ultimo film 1001
ultimo gara 821
ultimo generazione 893
ultimo giornata 1434
ultimo giorno 2436
ultimo giro 547
ultimo lavoro 707
ultimo libro 589
ultimo mese 1450
ultimo minuto 762
ultimo momento 1008
ultimo opera 798
ultimo ora 513
ultimo parola 965
ultimo parte 1134
ultimo partita 863
ultimo periodo 1200
ultimo piano 492
ultimo posto 903
ultimo puntata 1670
ultimo secolo 500
ultimo settimana 919
ultimo stagione 1765
ultimo tempo 1262
ultimo versione 1107
ultimo volta 2502
unico cosa 1931
unico differenza 559
unico eccezione 611
unico figlio 1210
unico modo 1675
unico navata 781
unico persona 643

unico punto 501	vecchio amico 1410
unico soluzione 630	vero identità 876
vario genere 1056	vero natura 586
vario gruppo 4402	vero nome 1980
vario tipo 1282	vero problema 540
vasto area 1056	vice presidente 757
vasto gamma 809	
vasto territorio 520	

Di seguito viene inoltre riportato l'insieme dei dati prodotti dallo strumento computazionale per tutte le espressioni considerate. I nomi delle colonne sono precisati dalla seguente legenda:

ESPRESSIONE	Espressione considerata
FREQ	Frequenza di occorrenza (componenti lemmatizzati)
FREQ_INT	Numero di occorrenze dell'espressione interrotta
INTERROMP	Valore di I_{syn}^{int}
FREQ_REV	Numero di occorrenze dell'espressione con componenti invertiti
ORDINE INVERSO	Valore di I_{syn}^{ord}
MOD_SINTAG	Valore di I_{syn}
FREQ_SIN1	Numero di occorrenze delle espressioni con primo componente sostituito (valore non intero poiché normalizzato in base alla similarità semantica)
DISP1	Dispersione della sostituzione del primo componente
FREQ_SIN2	Numero di occorrenze delle espressioni con secondo componente sostituito (valore non intero poiché normalizzato in base alla similarità semantica)
DISP2	Dispersione della sostituzione del secondo componente
MOD_PARAD	Valore di I_{sub}
FORMA_PREVAL	Forma con il più alto numero di occorrenze
FREQ_P	Numero di occorrenze della forma prevalente
FORMA_PREVAL_2	Forma con il secondo numero di più alto di occorrenze
FREQ_P_2	Numero di occorrenze della seconda forma prevalente
MOD_FLESS	Valore di I_{infl}
ERR	Se contrassegnato da X, espressione scartata per errori dovuti al corpus

emissario boreale	7477	913	0.55	0	0.55	0	0.511756545	0	0	0	1.4343751359	0.0019165009	emissario boreale	7477	0	1	0.0013386881	0
pianta erbacea	7477	2	0.002670227	781	0	0.5117647059	0	0	0	0	0.7171875019	0	pianta erbacea	7477	0	262	0.3574297189	0
pianta erbacea	7477	78	0.0945454545	0	0.0945454545	0	0.0945454545	0	0	0	0	0.0009591696	pianta erbacea	7477	0	162	0.2275799746	0
regione italiana	7477	78	0.0945454545	0	0.0945454545	1333	1840684096	0.0068647795	0	0	0.166698352	0.6412105068	regioni italiane	7477	0	165	0.2208835341	0
problema tecnico	7477	12	0.0158102767	0	0.0158102767	72	12462316764	0.2590314339	0	0	20.877264445	0.1103227834	problemi tecnici	7477	0	582	0.2208835341	0
produzione cinematografica	746	6	0.007987234	0	0.007987234	243	6544246014	0.490902271	0	0	3.9324134214	0.2491847662	produzione cinematografica	557	0	178	0.2533512064	0
caratteristica peculiare	744	89	0.0068427371	122	0.0068427371	73	0696790254	0.4850593231	0	0	683.3806237489	0.0337757162	caratteristica peculiare	477	0	302	0.4059139785	0
prodotto agricolo	744	3	0.0040160643	1	0.0040160643	49	7259518547	0.4135101496	0	0	1.6608056639	0.567704761	prodotto agricolo	697	0	45	0.063172043	0
pianta ciclabile	744	3	0.0040160643	0	0.0040160643	44	9387151772	0.0547319349	0	0	0.4621988391	0.3449035039	pianta ciclabile	744	0	369	0.501344086	0
rivera destra	744	44	0.0550905133	0	0.0550905133	391	2483482514	0.0147766753	0	0	0	0.022304127	letteratura italiana	741	0	1	0.00269179	0
letteratura italiana	743	44	0.0550905133	0	0.0550905133	1.6609011137	0	0.0147766753	0	0	0	0.022304127	letteratura italiana	741	0	1	0.00269179	0
favata laterale	743	13	0.017957672	0	0.017957672	0	0.017957672	0	0	0	2.0761264172	0	favata laterale	611	0	129	0.1778581427	0
scienza naturale	742	4	0.0035013303	0	0.0035013303	6.0950455195	1.0389207862	0	0	0	12.5499701975	0.1205822735	scienza naturale	693	0	41	0.003423181	0
tema principale	742	49	0.062052165	64	0.062052165	300.4217791384	0.19596000432	105.2641004267	0	0	103.656810761	0.116990486	tema principale	520	0	221	0.296245614	0
tempo indeterminato	741	5	0.0067024129	0	0.0067024129	8.9982431302	0.2175664451	0	0	0	0.1286989077	0.1319689041	tempo indeterminato	73	0	5	0.0094466937	0
diritto canonico	740	11	0.0146471372	0	0.0146471372	28.7378707977	0.21892618237	0	0	0	0.1970849163	0.0633126457	diritto canonico	73	0	1	0.0013513514	0
stile libero	740	10	0.0133333333	0	0.0133333333	123	0861498218	0.0620695322	0	0	11.5296255288	0.9175559018	stile libero	739	0	1	0.0013513514	0
popolazione civile	740	4	0.0053763441	0	0.0053763441	8.0557965703	0.5417381422	0	0	0	3.1869943826	0.0149655891	popolazione civile	655	0	84	0.1148648849	0
televisione digitale	740	0	0	0	0	42.1450301854	0.0569147149	0	0	0	0	0.0538839072	televisione digitale	732	0	8	0.108108108	0
faccia filosofica	739	0	0	0	0	0	0	0	0	0	0	0	faccia filosofica	739	0	0	0	0
tempo successivo	737	194	0.2043780881	112	0.1319199058	0.29339844679	31.878205019669	0.0035848319	0	0	81.6649854579	0.0953129487	tempo successivo	600	0	136	0.1858887381	0
vita reale	737	17	0.0256464191	8	0.0107388255	0.0328008398	172.5500992115	0.1453566428	24.3437294588	0	750.028029852	0.2104832582	vita reale	733	0	2	0.0054274084	0
architettura civile	737	4	0.005981107	0	0.005981107	260.9363675812	0.0419306981	39.1126625539	0.0869788724	0	29.92610044	0.2892610044	architettura civile	569	0	168	0.2279511533	0
impiego operativo	737	2	0.007063599	0	0.007063599	25.7995135189	0.8023324733	28.3213541857	0.0313365604	0	0.0660019395	0.0660019395	impiego operativo	737	0	0	0	0
comunità internazionale	736	120	0.1401869159	0	0.1401869159	1123.261641258	0.19124654231	37.5033032706	0.1924654231	0	0.6119780797	0.6119780797	comunità internazionale	730	0	6	0.0095108096	0
messe estivo	736	2	0.0027100271	0	0.0027100271	0	0	0	0.031175043	0	0.2205735748	0.2205735748	messe estivo	729	0	7	0.0013605442	0
mondo occidentale	735	27	0.0354330709	0	0.0354330709	208.001130364	0.0420107285	0	0	0	53.7051261751	0.0680065644	mondo occidentale	729	0	1	0.0013605442	0
beni culturali	735	27	0.0354330709	0	0.0354330709	233.2849838257	0.0339974213	0	0	0	1.6496107407	0.2422169531	beni culturali	699	0	35	0.4889795918	0
mondo esterno	735	3	0.0040650407	0	0.0040650407	445.3147739438	0.0325325056	51.129690098	0.3502792931	0	0.4031399489	0.4031399489	mondo esterno	730	0	5	0.0068072711	0

D

Dati sul pattern NPN

Si riporta di seguito la lista, in ordine alfabetico, delle 500 espressioni relative al pattern NPN estratte in PAISÀ e i dati relativi al comportamento variazionale di ognuna delle espressioni.

accordo di pace 366
aereo da caccia 400
aereo da trasporto 393
aereo di linea 463
agente di polizia 562
agenzia di stampa 367
albero a camma 546
album d' esordio 534
album da solista 389
album di debutto 1795
album di studio 369
album in studio 1056
amico d' infanzia 612
amico di famiglia 362
angolo di posizione 641
anime per titolo 930
anno d' età 488
anno di assenza 558
anno di attività 1095
anno di carcere 976
anno di carriera 767
anno di distanza 857
anno di età 1735
anno di guerra 667
anno di lavoro 642
anno di piombo 490
anno di prigionia 359
anno di reclusione 699

anno di regno 597
anno di servizio 422
anno di storia 524
anno di studio 538
anno di vita 2304
arco di tempo 680
arma da fuoco 1734
arma di distruzione 511
articolo con tag 422
asse di rotazione 399
atomo di carbonio 523
attacco da parte 520
attività di ricerca 706
atto di violenza 370
ausilio di strumento 765
azienda di trasporto 388
base di partenza 393
battuta di caccia 435
bocca da fuoco 374
borsa di studio 1605
cabina di pilotaggio 387
cacciatore di taglia 495
cadenza di tiro 553
calcio di rigore 954
camera da letto 654
campagna di scavo 398
campionato di calcio 1659
campionato di serie 885

- campione in carica 635
campo di battaglia 2137
campo di concentramento 2358
campo di gioco 768
campo di prigionia 440
campo di sterminio 478
canna da zucchero 499
capacità di carico 640
capo di stato 1053
capoluogo di provincia 646
carica di presidente 514
carrello d' atterraggio 491
carriera di allenatore 362
carriera di attore 611
carta di credito 950
casa di cura 354
casa di moda 373
casa di produzione 1474
casa di riposo 479
caso di necessità 567
causa di problema 725
cavallo di battaglia 599
centinaio di metro 1091
centinaio di migliaio 1191
centinaio di persona 626
centro di ricerca 834
cerimonia di premiazione 441
chiave di lettura 465
chilometro di distanza 520
ciclismo su strada 434
ciclo di affresco 609
ciclo di vita 388
cilindro in linea 950
classifica di vendita 508
codice di procedura 358
colpo di grazia 537
colpo di mano 411
colpo di pistola 998
colpo di scena 1388
colpo di stato 2183
colpo di testa 427
comandante in capo 768
community di tvblog 4234
compagno di classe 745
compagno di scuola 815
compagno di squadra 2733
compagno di viaggio 453
compositore di musica 398
comunicazione di massa 390
condanna a morte 842
condizione di lavoro 418
condizione di salute 873
condizione di vita 1375
conflitto d' interesse 436
conflitto di interesse 550
consiglio di amministrazione 1119
contratto di lavoro 484
controllo da parte 2977
corpo a corpo 363
corpo di fabbrica 571
corpo di spedizione 435
corso d' acqua 3258
corso di formazione 386
corso di laurea 850
corso di studio 458
costo di produzione 836
crimine di guerra 649
crisi di credibilità 361
critico d' arte 520
dato di fatto 475
datore di lavoro 1711
decina di anno 494
decina di metro 525
decina di migliaio 1206
densità di popolazione 943
dichiarazione di guerra 525
direttore d' orchestra 1235
diritto d' autore 1717
diritto di proprietà 426
disco d' oro 1005
disco di platino 1209
disegno di legge 1236
distretto di polizia 639
distruzione di massa 538
dozzina d' anno 376
edificio di culto 1008

-
- epidemia di pesta 460
episodio in onda 391
esempio di architettura 510
euro a persona 537
famiglia di origine 694
fascia d'età 794
fase a girone 499
fase di sviluppo 677
fatto di cronaca 454
fattore di rischio 590
febbre d'amore 519
ferro di cavallo 443
fibra di carbonio 707
figlio d'arte 730
film d'animazione 1369
film d'avventura 407
film d'azione 852
film di fantascienza 819
fin di vita 586
finale di stagione 299
foglio di carta 451
fonte di energia 630
fonte di ispirazione 564
forma d'arte 411
forma di governo 732
forma di vita 1243
forza di gravità 555
forza di polizia 1004
funzione d'onda 517
fuoco d'artificio 547
galleria d'arte 566
galleria di immagine 501
gara di andata 407
gas di scarico 438
gioco d'azzardo 752
gioco da tavolo 543
gioco di carta 663
gioco di parola 920
gioco di ruolo 1399
giornale di carta 1031
giornata di campionato 515
giorno d'oggi 1555
giorno di festa 372
girone d'andata 356
girone di andata 502
girone di ritorno 686
giuramento di fedeltà 372
gruppo di amico 520
gruppo di giovane 484
gruppo di lavoro 908
gruppo di persona 1297
gruppo di ragazzo 411
gruppo di ricerca 419
guerra d'indipendenza 902
guerra di indipendenza 1205
guerra di secessione 673
guerra di successione 924
immissione di nome 480
impianto di risalita 414
incidente d'auto 427
intervallo di tempo 637
intervento di restauro 454
istituto di credito 612
istituto di ricerca 423
lasso di tempo 908
lavorio di costruzione 794
lavorio di restauro 994
lavorio di ristrutturazione 679
libertà di espressione 649
libertà di ricerca 399
libertà di scelta 375
libertà di stampa 758
libro di testo 439
limite di velocità 382
linea di confine 401
linea di principio 398
linea di successione 514
linguaggio di programmazione 1186
locomotiva a vapore 433
lotta di classe 612
lunghezza d'onda 1939
luogo d'interesse 2968
luogo di culto 1862
luogo di interesse 1655
luogo di lavoro 565
luogo di nascita 569

macchina da presa 756	miliardo di dollaro 2506
maestro di cappella 561	miliardo di euro 2132
mal di testa 414	miliardo di lira 1039
mamma per amico 356	milione di abitante 951
mananza di fondo 446	milione di anno 4630
mano a mano 549	milione di copia 2308
marchio di fabbrica 455	milione di dollaro 6120
medaglia d' argento 2006	milione di euro 5413
medaglia d' oro 4785	milione di lira 806
medaglia di bronzo 1510	milione di persona 2728
media di rete 366	milione di spettatore 655
medico in famiglia 4919	milione di sterlina 792
mese di agosto 841	milione di telespettatore 1481
mese di aprile 675	milione di tonnellata 702
mese di dicembre 535	minuto per minuto 547
mese di febbraio 426	misura di sicurezza 481
mese di gennaio 480	modalità di gioco 906
mese di giugno 792	moneta d' argento 427
mese di luglio 938	moneta d' oro 575
mese di maggio 853	motivo di salute 502
mese di marzo 630	motivo di sicurezza 428
mese di novembre 568	motore a benzina 428
mese di ottobre 838	motore a combustione 387
mese di settembre 755	motore di ricerca 1183
mese di vita 498	muro di cinta 466
messa a punto 702	museo di storia 394
messa in onda 1743	musica da camera 469
messa in scena 899	nave da battaglia 1101
meta di pellegrinaggio 371	nave da guerra 1115
metro d' altezza 411	nazionale di calcio 446
metro di altezza 1043	nazionale di pallavolo 729
metro di altitudine 451	nome d' arte 1033
metro di diametro 400	nome di battaglia 390
metro di distanza 479	nome di battesimo 376
metro di lunghezza 998	nome in codice 1411
metro di profondità 697	not a valid 3528
metro di quota 387	numero di abitante 729
mezzo di comunicazione 1510	numero di lettore 384
mezzo di produzione 370	numero di persona 828
mezzo di trasporto 1536	numero di telefono 523
migliaio di anno 719	obiezione di coscienza 385
migliaio di persona 1507	oggetto di studio 706
miliardo di anno 719	olio d' oliva 503

-
- olio di oliva 369
olio su tela 919
ombra di dubbio 461
onda d' urto 465
opera d' arte 3367
opera di narrativa 573
ora di lavoro 352
ora di volo 474
ordine di grandezza 518
ordine di tempo 375
ottavo di finale 806
padrone di casa 1810
paese d' origine 525
paese in via 640
paio d' anno 543
paio di anno 474
paio di giorno 409
paio di mese 361
pala d' altare 733
pallone d' oro 385
parete di fondo 507
parola d' ordine 608
partita di andata 441
partita di calcio 642
partita di campionato 560
pazzo per amore 658
pena di morte 1505
pericolo di estinzione 397
periodo d' oro 488
periodo di crisi 603
periodo di tempo 1891
periodo di visibilità 1272
permesso di soggiorno 509
pesce d' acqua 431
piano di coda 389
pilota di formula 410
pò di tempo 1242
polvere da sparo 414
ponte di volo 394
popolazione in fascia 516
porta d' ingresso 519
porta di accesso 391
portale d' ingresso 540
portata di mano 440
porzione di territorio 514
posizione in classifica 502
posto a sedere 513
posto da titolare 366
posto di blocco 428
posto di lavoro 2129
posto in classifica 895
potenza di fuoco 672
presa d' aria 627
presa di coscienza 487
presa di posizione 732
prigioniero di guerra 741
principio di funzionamento 361
problema di salute 1007
processo di produzione 394
produzione di energia 838
produzione di serie 399
progetto di legge 459
progetto di ricerca 405
promozione in serie 518
proposta di legge 926
puntata di sabato 730
punto di contatto 512
punto di forza 1441
punto di incontro 374
punto di morte 554
punto di partenza 1927
punto di riferimento 3307
punto di svolta 516
punto di vantaggio 631
punto di vista 24159
quantità di energia 498
quantità di moto 494
quota di mercato 494
raccolta di poesia 522
raccolta di racconto 413
ragazzo di nome 395
raggio d' azione 857
rango di arcidiocesi 473
rapporto di amicizia 549
rapporto di compressione 864
rapporto di forza 373

- rapporto di lavoro 638
reazione a catena 433
reggimento di fanteria 361
richiesta di chiarimento 2804
risultato di rilievo 384
romanzo di fantascienza 805
ruolo da protagonista 430
ruolo di protagonista 389
salto di qualità 739
sci di fondo 492
scontro a fuoco 428
scopo di lucro 481
scuola di pensiero 430
secondo d' arco 748
sedia a rotella 683
segretario di stato 392
seguito di controllo 483
senso di colpa 1042
senso di marcia 560
serie a fumetto 385
serie di concerto 595
serie di evento 464
serie di film 402
serie di videogioco 531
servizio di trasporto 445
sistema d' arma 598
sistema di comunicazione 432
sistema di controllo 1208
sistema di gestione 407
sistema di guida 525
sistema di navigazione 425
sistema di riferimento 782
sistema di scrittura 397
sistema di sicurezza 469
sistema di trasporto 484
sito di interesse 441
società per azione 652
soluzione di continuità 535
somma di denaro 961
specchio d' acqua 534
specie di uccello 1065
squadra di calcio 1962
stato d' animo 773
stato di abbandono 445
stato di conservazione 1094
stato di cosa 363
stato di diritto 415
stato di salute 754
stella di classe 2809
stella di magnitudine 473
stile di vita 3121
storia a fumetto 584
storia d' amore 2290
strumento di misura 438
studio di registrazione 846
successo di critica 423
successo di pubblico 1241
successo di vendita 562
tasso d' interesse 362
tasso di cambio 429
tasso di crescita 446
tasso di interesse 583
tempesta d' amore 1486
tempo di guerra 567
tempo di pace 457
tempo di percorrenza 396
tenore di vita 389
terreno di gioco 442
territorio a vantaggio 893
tesi di laurea 463
testa di ponte 442
testa di serie 389
tipo di arma 357
tipo di dato 410
titolo di campione 1070
titolo di coda 585
titolo di conte 599
titolo di duca 407
titolo di testa 398
traffico di droga 425
unità di misura 1340
uomo d' affare 1099
uscita di scena 437
uso di droga 397
valvola per cilindro 556
velocità di rotazione 421

vettura di formula 637
 via d' uscita 489
 via di comunicazione 1302
 via di fuga 629
 via di mezzo 495
 via di sviluppo 1217
 viaggio di ritorno 603
 vicino di casa 354

vincitore di medaglia 1308
 visibilità ad occhio 756
 volo di linea 410
 volta a botte 578
 volta a crociera 440
 volta in volta 452

Di seguito viene inoltre riportato l'insieme dei dati prodotti dallo strumento computazionale per tutte le espressioni considerate. I nomi delle colonne sono precisati dalla seguente legenda:

ESPRESSIONE	Espressione considerata
FREQ	Frequenza di occorrenza (componenti lemmatizzati)
FREQ_INT1	Numero di occorrenze dell'espressione interrotta tra primo e secondo componente
FREQ_INT2	Numero di occorrenze dell'espressione interrotta tra secondo e terzo componente
MOD_SINTAG	Valore di I_{syn}
FREQ_SIN1	Numero di occorrenze delle espressioni con primo componente sostituito (valore non intero poiché normalizzato in base alla similarità semantica)
DISP1	Dispersione della sostituzione del primo componente
FREQ_SIN2	Numero di occorrenze delle espressioni con secondo componente pieno sostituito (valore non intero poiché normalizzato in base alla similarità semantica)
DISP2	Dispersione della sostituzione del secondo componente
MOD_PARAD	Valore di I_{sub}
FORMA_PREVAL	Forma con il più alto numero di occorrenze
FREQ_P	Numero di occorrenze della forma prevalente
FORMA_PREVAL_2	Forma con il secondo numero di più alto di occorrenze
FREQ_P_2	Numero di occorrenze della seconda forma prevalente
MOD_FLESS	Valore di I_{infl}
ERR	Se contrassegnato da X, espressione scartata per errori dovuti al corpus

[illegible]

367	0	0	0	0	11.5223139567	0.1378921355	0	0.0289125942	cabina di pilotaggio	376	cabine di pilotaggio	11	0.028423726
367	1	0	0.0025773196	0.015251909	0	501.51.76957935	0	0.0126609639	383metro di quota	384	metro di quota	4	0.0103359173
367	16	1	0.0420792079	0.1860089207	0	0	0	0.0004804123	motore a combustione	194	motore a combustione	193	0.4987080103
366	9	12	0.0515970516	105.9832878908	0.2085438511	64.447589472	0	0.4953089041	263corso di formazione	263	corso di formazione	122	0.3186528497
365	2	0	0.0051679587	0	0	0	0	0	obiezione di coscienza	379	obiezioni di coscienza	6	0.0158844156
365	0	0	0	0	0	0	0	0	pallore d' oro	332	pallore d' oro	53	0.1376623377
365	14	0	0.0350877193	0	0	0.9264237425	0	0.0024005191	serie a funetto	384	serie a funetto	1	0.0025974026
364	0	0	0.0327455919	2.8608370833	1.7298516228	0.0010191443	0	0.007397618	numero di lettere	195	numero di lettere	158	0.4114583333
364	12	66	0.1688311668	16.8460073322	0.3659958047	2.3240113357	0.3391077072	0.0475787248	risultato di rilievo	226	risultato di rilievo	187	0.4895287958
362	38	2	0.0947867299	85.0539733224	0.2076065185	0	0	0.1821073756	limiti di velocità	195	limite di velocità	2	0.0053191489
376	0	0	0	0	0	0	0	0	dozzina d' anno	374	dozzine d' anni	2	0.003793935
376	2	2	0.0105263158	0.142908652	0	0	0	0.000379935	nome di battesimo	350	nomi di battesimo	26	0.0691489362
375	5	0	0.0131578947	207.0934025934	0.0379290593	8.5043924741	0.4951968899	0.3850501185	libertà di scelta	369	libertà di scelta	6	0.016
375	1	0	0.0026595745	1.9182161203	0.6929389718	12.198368267	0.2965427002	0.0062785472	ordine di tempo	374	ordini di tempo	1	0.0026666667
374	0	0	0	0	0	0.2369419858	0	0.0006331336	bocca da fuoco	197	bocche da fuoco	177	0.4732620321
374	2	0	0.0053191448	269.8861219812	0.0181775584	194.929900492	0.0577506081	0.5541334572	punto di incontro	347	punti di incontro	27	0.0721925134
373	2	1	0.0079787234	0	0	0.0692391893	1.4649288038	0.0001855934	casa di moda	213	case di moda	156	0.4289544236
373	3	0	0.0026737968	2.558242573	0.236615708	2.5920607087	2.5840989946	0.0136197254	rapporto di forza	270	rapporto di forza	64	0.2761394102
372	3	8	0.0287206266	63.4353646232	0.736615708	99.8984215666	1.062954748	0.3051064413	giorno di festa	211	giorno di festa	159	0.4327956989
372	5	2	0.018469657	12.5294075183	1.0286334867	0.4651352243	1.8427888725	0.0337500277	giuramento di fedeltà	359	giuramenti di fedeltà	17	0.0456989247
371	13	88	0.2139830508	0	0	11.8765423993	0	0.0310192479	metà di pellegrinaggio	191	metà di pellegrinaggio	154	0.4851752022
370	7	14	0.0537084399	31.8771281996	0.6183020362	26.9681493301	0.7938458681	0.1372178679	atto di violenza	272	atto di violenza	98	0.2648646489
370	2	0	0.0053763441	288.0708195887	0.059507047	30.1813758024	0.2747972033	0.4624063643	mezzo di produzione	350	mezzo di produzione	19	0.0540540541
369	160	0	0.3024574669	0	0	0	0	0	olio di oliva	361	oli di oliva	6	0.0216802168
367	5	0	0.0134408602	0.2061714074	0	0	0	0.0005614595	agenzia di stampa	210	agenzie di stampa	157	0.4277929155
366	7	0	0.018766756	10.2929000554	0.8037622743	7.0214710825	0.287437191	0.0451703101	accordo di pace	211	accordi di pace	155	0.4234972878
366	7	3	0.0265957447	0	0	5.1696629697	0.2373776474	0.0139280852	media di rete	359	medie di rete	7	0.0191256631
366	5	1	0.0161290323	0.2060959578	0	0.0850166286	0	0.0007947574	posto da titolare	365	posto da titolari	1	0.0027322404
363	0	0	0	0	0	1.4142135624	0	0	corpo a corpo	359	corpi a corpo	4	0.0110192837
363	4	1	0.0135869565	0.0053494431	0.667261314	75.3411044843	0.140242855	0.17564909	stato di cosa	360	stati di cose	2	0.0082644628
362	1	2	0.0082191781	0.5633212713	1.9687262491	1.8306112501	0.9585824814	0.0055996278	amico di famiglia	241	amici di famiglia	67	0.3342541436
362	2	0	0.0054945055	0.3032882669	1.6670061074	0	0	0.0008371115	carriera di allenatore	359	carriere di allenatrice	3	0.0082872928
362	4	0	0.0109289617	0	0	12.24730555	0.4815130945	0.0327251669	lasso d' interesse	202	lassi d' interesse	159	0.4419899503
361	0	0	0	0	0	0	0	0	crisi di credibilità	361	crisi di credibilità	0	X
361	0	0	0	0	0	0	0	0	paio di mese	360	paio di mese	1	0.0027700831
361	4	0	0.0109589041	0	0	0	0	0	principio di funzionamento	361	principio di funzionamento	0	0
361	5	0	0.0136612022	0	0	0	0	0	reggimento di fanteria	250	leggimenti di fanteria	111	0.3074792244
359	0	0	0	0.3848025758	0	1068.218302714	0.005666775	0.7485295467	anno di prigione	329	anno di prigione	30	0.0835654596
358	4	0	0.0110497238	5.0260500489	0.2352258569	3.328127076	0.406803282	0.0228035536	codice di procedura	350	codici di procedura	8	0.0223463687
357	1	3	0.0110803324	48.4048300781	0.6151996671	4.489108084	2.2775354829	0.1290429883	tipo di arma	155	tipi di armi	119	0.0558263305
356	0	0	0	0	0	0	0	0	giorno d' andata	354	giorni d' andata	2	0.0056179775
356	0	0	0	0	0	0	0	0	manina per amico	354	manina per amico	2	0.0056179775
354	0	0	0	0	0	2.7811915564	2.0281946954	0.0077952303	casa di cura	282	case di cura	70	0.2033898305
354	1	0	0.0028169014	0	0	0	0	0	vicino di casa	303	vichi di casa	28	0.1440677966
352	4	10	0.0382513661	125.4270058037	0.0616554235	30.2741575551	1.6041016557	0.3066787603	ora di lavoro	322	ora di lavoro	30	0.0852272727

E

Dati sul pattern NPdN

Si riporta di seguito la lista, in ordine alfabetico, delle 500 espressioni relative al pattern NPdN estratte in PAISÀ e i dati relativi al comportamento variazionale di ognuna delle espressioni.

24,63- su target 2013	ascesa al trono 645
abitante di città 394	aspetto di vita 466
abitante di luogo 346	attacco al suolo 503
abitante di paese 314	attenzione di pubblico 304
abitante di provincia 4270	aumento di popolazione 304
abitante di villaggio 364	aumento di prezzo 438
abitante di zona 401	autore di libro 406
accettazione di regola 485	autore di testo 334
acqua di fiume 530	azienda di provincia 808
acqua di lago 301	base al media 345
addetto al lavoro 1087	battaglia di guerra 363
album di band 657	biografia di personaggio 883
album di cantante 357	bocca al lupo 368
album di gruppo 669	bordo di nave 439
allevamento di bestiame 375	caduta di governo 294
alwaysurbi al ora 324	caduta di impero 302
ambito di progetto 313	caduta di muro 349
ambito di programma 292	caduta di regime 339
ambito di sito 483	campionato di mondo 2841
anno di guerra 384	campione di mondo 3528
anno di morte 303	canzone di album 359
anteprima di commento 367	capitale di regno 384
appartenente al famiglia 1045	capitale di stato 329
appartenente al genere 328	capitano di squadra 389
applicazione di legge 331	capitolo di saga 485
arcidiocesi al termine 422	capitolo di serie 371
artista di calibro 357	capo di esercito 382
ascesa al potere 400	capo di governo 896

- capo di polizia 400
capo di stato 472
capoluogo di provincia 347
casa di libertà 334
castello di giornale 337
castello di provincia 381
causa di abuso 512
causa di difficoltà 324
causa di fatto 491
causa di guerra 523
causa di mancanza 590
causa di morte 580
causa di presenza 430
cavaliere di zodiaco 587
centro di attenzione 459
centro di città 1219
centro di paese 596
chiesa di provincia 498
cima di monte 295
codice di strada 384
colpo al minuto 323
comando di generale 466
commento da contenuto 480
commento di lettore 4945
commento in ambito 480
complessità di compito 350
componenti di gruppo 612
comuni di comunità 381
comuni di distretto 397
comuni di provincia 5829
comuni di regione 427
coniuge di re 317
conoscenza di fatto 347
consapevolezza di fatto 318
conseguenza di fallimento 301
consiglio di ministro 392
conto di fatto 308
controllo di territorio 439
copertina di album 414
coppa di mondo 448
corso di anno 5957
corso di episodio 442
corso di fiume 1130
corso di guerra 843
corso di secolo 2443
corso di serie 366
corso di stagione 615
corso di storia 939
corso di tempo 1154
cosa di genere 692
costruzione di chiesa 509
crescita di popolazione 307
decennio di secolo 299
decreto di presidente 443
difesa di città 378
difesa di diritto 457
dinamica di idea 332
diocesi al termine 1134
direttore di fotografia 372
direttore di istituto 437
diritto di donna 343
diritto di lavoro 297
diritto di uomo 720
economia di conoscenza 885
economia di paese 303
ecosistema di informazione 678
edizione di campionato 492
effetto di inquinamento 602
effetto di moltiplicazione 292
elenco di patrimonio 293
episodio di serie 1514
erede al trono 1102
erezione di diocesi 727
età di bronzo 1190
età di ferro 681
età di oro 328
età di pietra 309
etimologia di nome 305
euro al anno 176
euro al mese 432
facciata di chiesa 458
fallimento di sistema 308
fatto di giorno 300
figlio di conte 677
figlio di duca 462
figlio di imperatore 378

-
- figlio di re 1533
figlio di uomo 376
film di serie 378
fin di conte 391
finale di campionato 379
fine di anno 9956
fine di campionato 408
fine di conflitto 760
fine di estate 336
fine di film 363
fine di guerra 3211
fine di mese 498
fine di mondo 623
fine di periodo 384
fine di secolo 1189
fine di serie 318
fine di stagione 1110
fine di storia 309
fiore al occhiello 384
fisica di particella 303
fiume di provincia 664
foce di fiume 791
fondazione di città 336
forza di ordine 3285
frazione di comune 2153
frazione di provincia 1750
gara di campionato 288
gara di stagione 289
giorno di settimana 325
giro al minuto 422
giro di mondo 534
grazie al aiuto 873
grazie al collaborazione 245
grazie al declinazione 1080
grazie al fatto 387
grazie al intervento 543
grazie al opera 271
grazie al posizione 1315
grazie al presenza 577
grazie al successo 226
guardia di corpo 1121
inizio di anno 5941
inizio di campionato 327
inizio di carriera 389
inizio di estate 316
inizio di film 311
inizio di guerra 786
inizio di inverno 415
inizio di lavoro 366
inizio di secolo 867
inizio di stagione 870
inizio di storia 301
interno di area 338
interno di casa 312
interno di chiesa 1125
interno di città 437
interno di comunità 322
interno di corpo 298
interno di edificio 416
interno di famiglia 301
interno di gruppo 1268
interno di mura 385
interno di parco 378
interno di struttura 323
interno di territorio 427
ironia di sorte 401
largo di costa 614
lavaggio di cervello 344
leader di gruppo 362
leader di partito 351
lettera di alfabeto 538
limite di visibilità 759
lista di ospite 430
livello di acqua 350
livello di mare 2781
luce di sole 787
lunedì al venerdì 518
macchina di tempo 397
maggioranza di caso 416
maggioranza di popolazione 592
massimo di voto 343
membro di band 825
membro di comitato 467
membro di commissione 373
membro di consiglio 550
membro di equipaggio 776

membro di famiglia 1923	occasione di festa 441
membro di governo 324	oggetto di catalogo 414
membro di gruppo 1535	opera di architetto 501
membro di parlamento 326	opera di artista 397
membro di partito 319	opera di pittore 353
membro di squadra 374	opera di scultore 426
mercato di lavoro 720	ora al giorno 497
mese di anno 395	ora di giorno 299
metà di anno 6961	ordine di giorno 1149
metà di secolo 1744	origine di nome 1569
metro su livello 1002	origine di termine 352
miglioramento di condizione 325	osservatore di emisfero 831
ministro di difesa 314	paese di mondo 764
ministro di giustizia 316	parte di abitante 332
ministro di guerra 438	parte di anno 600
ministro di interno 471	parte di area 488
moltiplicazione di messaggio 296	parte di autorità 411
momento di lancio 289	parte di caso 1666
momento di morte 306	parte di cast 624
mondo di calcio 401	parte di città 832
mondo di cinema 768	parte di colonna 344
mondo di lavoro 686	parte di complesso 330
mondo di moda 303	parte di comune 419
mondo di musica 625	parte di comunità 427
mondo di spettacolo 1429	parte di corpo 1149
morte di figlio 425	parte di critica 522
morte di fratello 598	parte di edificio 359
morte di madre 654	parte di esercito 674
morte di marito 510	parte di famiglia 672
morte di moglie 441	parte di film 409
morte di padre 2425	parte di forza 745
morte di re 334	parte di governo 942
mura di città 578	parte di gruppo 1973
nome di band 306	parte di materiale 291
nome di città 763	parte di mondo 2354
nome di genere 378	parte di movimento 318
nome di gruppo 535	parte di opera 374
nome di paese 405	parte di paese 602
nome di personaggio 346	parte di patrimonio 326
nord di città 454	parte di persona 480
nord di paese 286	parte di popolazione 1660
notizia di morte 469	parte di progetto 394
occasione di elezione 398	parte di programma 403

-
- parte di provincia 1235
parte di pubblico 687
parte di regione 1839
parte di regno 346
parte di serie 350
parte di sistema 694
parte di società 425
parte di specie 381
parte di squadra 682
parte di stagione 465
parte di stato 650
parte di storia 502
parte di struttura 325
parte di studioso 333
parte di tempo 448
parte di territorio 1729
parte di testo 342
parte di truppa 549
parte di uomo 412
parte di utente 324
patrimonio di umanità 1211
patrono di città 398
pendice di monte 440
periodo di anno 635
periodo di guerra 299
personaggio di fumetto 1483
personaggio di mitologia 769
personaggio di serie 1374
personaggio di universo 409
piede di monte 334
popolo di libertà 958
porta di città 491
posizione di classifica 508
posizione di stella 1327
posizione in classifica 506
posto a sole 1235
posto al sole 929
posto di classifica 1001
posto in classifica 520
potenza di motore 390
prefetto di pretorio 380
presa di potere 341
presidente di associazione 272
presidente di commissione 390
presidente di consiglio 1166
presidente di provincia 622
presidente di regione 308
presidente di repubblica 494
pressi di circolo 336
pressi di città 433
pretendente al trono 321
progetto di architetto 808
proprietà di famiglia 670
prossimità di circolo 477
prossimità di equatore 615
protagonista di film 614
protagonista di serie 545
pubblicazione di album 389
qualità di prodotto 360
qualità di servizio 290
qualità di vita 925
racconto di innovazione 332
record di mondo 537
registrazione di album 353
regola di gioco 388
resa di conte 393
responsabilità di commento 480
responsabilità su contenuto 2743
resto di corpo 372
resto di gruppo 384
resto di mondo 1765
resto di paese 296
rete in club 747
rettore di università 274
ridefinizione di ruolo 296
rimozione di video 2746
risarcimento di danno 358
rispetto al anno 479
rispetto al media 309
rispetto al modello 323
rispetto al passato 480
rispetto al resto 445
rispetto al versione 492
rispetto di diritto 413
ritiro di truppa 318
riva di fiume 915

- riva di lago 344
romanzo di scrittore 449
ruolo di giornalismo 345
ruolo di protagonista 394
salone di automobile 305
salto in passato 384
scena di film 524
sceneggiatura di film 304
sciopero di fame 557
scopo di gioco 350
scoppio di guerra 988
sede di chiesa 1085
sede di governo 291
seguito al morte 548
seguito di morte 309
senso di articolo 519
senso di legge 534
senso di umorismo 303
senso di vita 299
significato di nome 538
signore di guerra 449
simbolo di città 409
sindaco di città 344
sintomo di difficoltà 338
società di gruppo 308
soluzione di equazione 326
soluzione di problema 397
specie di genere 828
sponda di fiume 432
stagione di pioggia 349
stagione di serie 701
stazione di metropolitana 324
stemma di famiglia 318
storia di arte 982
storia di calcio 384
storia di cinema 924
storia di città 391
storia di gruppo 711
storia di letteratura 320
storia di musica 857
storia di scienza 469
storia di umanità 294
strada di città 363
studente di scuola 295
successione al trono 403
sud di città 449
suddivisione di popolazione 552
supplemento al testata 3390
taglio di costo 392
tecnico di suono 316
tempo di guerra 331
teoria di complotto 353
teoria di insieme 406
teoria di numero 479
teoria di relatività 432
termine di anno 1690
termine di campionato 463
termine di conflitto 516
termine di guerra 562
termine di stagione 1822
terremoto in mezzo 337
territorio da diocesi 492
territorio di comune 990
territorio di provincia 392
testo di canzone 750
titolo di album 367
traccia di album 324
tratto da romanzo 304
tutela di diritto 341
ufficiale di esercito 398
unione di parola 433
università di studio 395
uscita di album 634
uscita di film 383
uso di arma 373
uso di forza 407
uso di termine 465
valle di fiume 902
valle di torrente 299
vantaggio di erezione 1338
velocità di luce 682
versione di fatto 308
vetta di classifica 521
via di centro 335
via di città 612
via di paese 379

via di ritorno 304	vita di uomo 311
viaggio in tempo 336	vittoria di campionato 400
video di canzone 318	volta al anno 382
vigile di fuoco 779	volta al giorno 400
vincitore di premio 540	volta in storia 658
violazione di diritto 664	zona di città 299
visione di mondo 643	
visualizzazione di video 2753	

Di seguito viene inoltre riportato l'insieme dei dati prodotti dallo strumento computazionale per tutte le espressioni considerate. I nomi delle colonne sono precisati dalla seguente legenda:

ESPRESSIONE	Espressione considerata
FREQ	Frequenza di occorrenza (componenti lemmatizzati)
FREQ_INT1	Numero di occorrenze dell'espressione interrotta tra primo e secondo componente
FREQ_INT2	Numero di occorrenze dell'espressione interrotta tra secondo e terzo componente
MOD_SINTAG	Valore di I_{syn}
FREQ_SIN1	Numero di occorrenze delle espressioni con primo componente sostituito (valore non intero poiché normalizzato in base alla similarità semantica)
DISP1	Dispersione della sostituzione del primo componente
FREQ_SIN2	Numero di occorrenze delle espressioni con secondo componente pieno sostituito (valore non intero poiché normalizzato in base alla similarità semantica)
DISP2	Dispersione della sostituzione del secondo componente
MOD_PARAD	Valore di I_{sub}
FORMA_PREVAL	Forma con il più alto numero di occorrenze
FREQ_P	Numero di occorrenze della forma prevalente
FORMA_PREVAL_2	Forma con il secondo numero di più alto di occorrenze
FREQ_P_2	Numero di occorrenze della seconda forma prevalente
MOD_FLESS	Valore di I_{infl}
ERR	Se contrassegnato da X, espressione scartata per errori dovuti al corpus

294	22	2	0.0754716981	0.2250762191	0	189.1314068895	0.0295608799	0.3917532706	storia dell' umanità	289	store dell' umanità	3	0.0170068027
299	7	664	0.6917525773	0	0	1.3021899986	1.0484159831	0.0043362654	decenni del secolo	190	decennio del secolo	109	0.364548495
299	76	14	0.2313624679	29.8351700726	0.1963518298	84.0235309717	0.0548740997	0.2757812809	ora di giorno	163	ora del giorno	136	0.4548494983
299	19	189	0.4102564103	486.9603558131	0.0350502963	18.8514486107	0.8179843732	0.6284845745	periodo della guerra	233	periodo delle guerre	60	0.220735786
274	2	9	0.0385964912	2.0717648789	0.6708191581	14.4359269223	0.495476505	0.0566235963	retore dell' università	242	rettori delle università	24	0.1167883212
299	24	21	0.1306139535	86.2252096634	0.5346026937	28.6734003298	0.2540910977	0.2776008596	sensò della vita	299	sensò delle vita	1	0.00334444816
299	9	4	0.0416666667	47.1600632014	0.31100538	25.76597973178	0.140781149	0.1960763374	valle dei torrente	266	valli dei torrenti	27	0.110367893
299	574	13	0.6625282167	910.6893008489	0.0100978447	4.9662273973	0.8700981453	0.7938396747	zona della città	175	zona della città	121	0.4147157191
298	1	91	0.2358974359	0.5835015506	1.996773894	574.00121047	0.0598280583	0.6584858801	interno dei corpi	282	interno dei corpi	11	0.0536912752
297	6	0	0.0198019802	19.1226724923	1.3589202747	13.2347760319	0.7484162052	0.0982441681	diritto del lavoro	276	diritti del lavoro	21	0.0707070707
296	0	0	0	0.4329433221	0	0	0	0.0014605101	moltiplicazione dei messaggi	296		0	X
286	7	4	0.037037037	7.963705693	0	14.3029105074	0.4084756846	0.07223168	nord del paese	283	nord dei paesi	3	0.0104895105
296	4	0	0.0133333333	0	0	47.6206700684	0.1724763336	0.1398585499	resto del paese	264	resto dei paesi	30	0.1081081081
296	0	2	0.0667114094	0	0	0.0020965744	1.4156976201	7.08E-006	ridefinizione dei ruoli	297	ridefinizione del ruolo	4	0.0195135135 X
295	14	13	0.0838509317	200.0397420645	0.0754197488	93.9090291525	0.0294842049	0.4991075381	cima del monte	235	cime dei monti	42	0.20333898305
295	10	27	0.1114457831	202.77628789	0.0419146965	70.2850730231	0.2156196931	0.4806716195	studenti della scuola	148	studenti della scuola	76	0.4983050847
294	7	55	0.1741573034	69.6466410664	0.5239562682	8.8219073844	0.453330523	0.2106716091	caduta del governo	292	caduta dei governi	2	0.0068027211
272	28	9	0.1197411003	30.884774666	0.6417206476	726.2297499028	0.162133061	0.7356951106	presidente dell' associazione	241	presidente della associazione	14	0.1139705882
293	0	4	0.0134680135	162.2317775867	0.6492344348	9.6317007663	0.1222482398	0.3697074224	elenco dei patrimoni	290	elenco del patrimonio	3	0.0102369078 X
292	0	24	0.0759493671	2.865663073	0.4223646161	164.3257867851	0.0101155407	0.3640998685	ambito del programma	274	ambito dei programmi	18	0.0616438356
292	1	0	0.0034129693	2.5534877114	0	4.8464091059	0.3036469749	0.0247157465	effetto della moltiplicazione	291	effetti della moltiplicazione	1	0.0034246575
291	3	22	0.0791139241	0	0	202.0866094205	0.1711694967	0.4098398866	parte del materiale	239	parte dei materiali	49	0.1786941581
291	13	7	0.0643086817	7.2656890337	0.4417206696	242.4474728604	0.2411045885	0.4618218669	sede del governo	282	sedì del governo	7	0.0309278351
290	6	3	0.0301003344	14.4856043944	0.3906531755	19.322323789	0.3308118744	0.1044070998	qualità del servizio	168	qualità dei servizi	122	0.4206896552
289	33	1	0.1057631579	6.7205271482	0.854657223	0.0488328365	0	0.0226872764	gara della stagione	175	gara della stagione	114	0.3944636678
289	0	10	0.0334448161	0.6910074865	2.160766464	10.9857511905	0.5978062183	0.0398349227	momento del lancio	289		0	0
288	29	8	0.1028037383	6.921471286	0.9261986592	0	0	0.0234688619	gara del campionato	160	gara del campionato	122	0.4444444444

F

Dati sul pattern NPV_{inf}

Si riporta di seguito la lista, in ordine alfabetico, delle 500 espressioni relative al pattern NPV_{inf} estratte in PAISÀ e i dati relativi al comportamento variazionale di ognuna delle espressioni.

accusa di essere 226	capacità di volare 103
associazione a delinquere 199	caso di dire 229
associazione per delinquere 135	caso di fare 102
attesa di vedere 112	compito di controllare 94
autorizzazione a procedere 138	compito di difendere 119
bisogno di andare 73	compito di fare 178
bisogno di avere 91	compito di gestire 83
bisogno di dire 63	compito di guidare 71
bisogno di essere 181	compito di organizzare 82
bisogno di fare 206	compito di portare 102
bisogno di sapere 67	compito di proteggere 141
bisogno di trovare 65	compito di rispondere 354
capacità di agire 86	consapevolezza di essere 74
capacità di assorbire 67	conto di avere 111
capacità di comunicare 66	conto di essere 160
capacità di creare 125	coraggio di dire 237
capacità di dare 67	coraggio di fare 123
capacità di essere 78	cosa da dire 81
capacità di fare 225	cosa da fare 321
capacità di gestire 73	decisione di abbandonare 68
capacità di leggere 71	decisione di fare 100
capacità di operare 66	decisione di lasciare 113
capacità di parlare 75	desiderio di avere 81
capacità di portare 83	desiderio di conoscere 65
capacità di produrre 124	desiderio di diventare 88
capacità di trasformare 74	desiderio di fare 98
capacità di trasportare 69	desiderio di vedere 74
capacità di vedere 69	diritto di avere 65

diritto di decidere 79
 diritto di dire 75
 diritto di eliminare 483
 diritto di essere 79
 diritto di fare 188
 diritto di manifestare 80
 diritto di nominare 74
 diritto di partecipare 145
 diritto di scegliere 100
 diritto di vivere 89
 dovere di fare 90
 fama di essere 118
 fatto di avere 310
 fatto di essere 780
 fatto di trovare 77
 fine di assicurare 73
 fine di creare 77
 fine di evitare 217
 fine di fare 110
 fine di garantire 129
 fine di ottenere 198
 fine di rendere 136
 fine di ridurre 90
 finta di essere 133
 fortuna di avere 78
 gioia di vivere 146
 gomma da masticare 71
 grado di accogliere 84
 grado di adattare 74
 grado di affrontare 215
 grado di agire 71
 grado di aiutare 106
 grado di aprire 78
 grado di assicurare 146
 grado di assorbire 147
 grado di assumere 82
 grado di attaccare 74
 grado di attivare 63
 grado di aumentare 112
 grado di avere 92
 grado di battere 72
 grado di bloccare 86
 grado di cambiare 83

grado di camminare 63
 grado di capire 200
 grado di catturare 71
 grado di colpire 100
 grado di combattere 123
 grado di competere 230
 grado di compiere 174
 grado di comprendere 143
 grado di comunicare 140
 grado di condurre 66
 grado di contenere 146
 grado di contrastare 136
 grado di controllare 310
 grado di coprire 82
 grado di costruire 107
 grado di creare 268
 grado di curare 70
 grado di dare 440
 grado di determinare 116
 grado di difendere 113
 grado di dimostrare 95
 grado di dire 125
 grado di distinguere 124
 grado di distruggere 159
 grado di effettuare 224
 grado di eliminare 66
 grado di emettere 94
 grado di erogare 579
 grado di eseguire 345
 grado di esercitare 78
 grado di esprimere 140
 grado di fare 1245
 grado di fermare 99
 grado di formare 71
 grado di fornire 636
 grado di funzionare 147
 grado di garantire 288
 grado di generare 311
 grado di gestire 410
 grado di identificare 73
 grado di individuare 118
 grado di indurre 67
 grado di influenzare 81

-
- grado di interpretare 71
grado di lanciare 153
grado di lavorare 100
grado di legare 111
grado di leggere 250
grado di manipolare 77
grado di mantenere 207
grado di mettere 199
grado di misurare 95
grado di modificare 109
grado di mostrare 106
grado di muovere 188
grado di offrire 345
grado di operare 205
grado di opporre 79
grado di ospitare 190
grado di ottenere 128
grado di pagare 115
grado di parlare 184
grado di passare 65
grado di penetrare 81
grado di percepire 100
grado di perforare 71
grado di portare 289
grado di prendere 125
grado di prevedere 65
grado di produrre 545
grado di proteggere 100
grado di provocare 91
grado di raccogliere 67
grado di raggiungere 451
grado di rappresentare 78
grado di realizzare 141
grado di reggere 129
grado di registrare 95
grado di rendere 178
grado di resistere 189
grado di ricevere 116
grado di riconoscere 219
grado di ricostruire 65
grado di ridurre 112
grado di rilevare 156
grado di riprodurre 247
grado di risolvere 168
grado di rispondere 145
grado di salvare 80
grado di sapere 122
grado di sconfiggere 98
grado di scrivere 64
grado di seguire 71
grado di sfruttare 154
grado di soddisfare 195
grado di sopportare 135
grado di sopravvivere 124
grado di sostenere 233
grado di sostituire 90
grado di sparare 102
grado di spiegare 132
grado di spingere 96
grado di spostare 73
grado di stabilire 103
grado di suonare 76
grado di superare 180
grado di supportare 117
grado di sviluppare 236
grado di svolgere 205
grado di tenere 188
grado di trasformare 196
grado di trasmettere 115
grado di trasportare 300
grado di trovare 118
grado di uccidere 129
grado di usare 171
grado di utilizzare 233
grado di valutare 88
grado di vedere 162
grado di vincere 107
grado di visualizzare 66
grado di volare 237
idea di avere 70
idea di costruire 92
idea di creare 199
idea di dare 72
idea di essere 78
idea di fare 282
idea di realizzare 123

idea di utilizzare 77	modo di affrontare 89
impressione di essere 145	modo di agire 158
incarico di formare 82	modo di apprezzare 69
intento di creare 104	modo di cantare 120
intento di fare 178	modo di combattere 103
intenzione di abbandonare 64	modo di comportare 79
intenzione di creare 95	modo di comunicare 85
intenzione di dare 77	modo di concepire 205
intenzione di fare 424	modo di conoscere 526
intenzione di lasciare 133	modo di dare 70
intenzione di portare 70	modo di dimostrare 66
intenzione di realizzare 73	modo di dire 783
intenzione di sposare 82	modo di entrare 72
intenzione di tornare 80	modo di esprimere 219
intenzione di uccidere 102	modo di essere 450
libertà di scegliere 65	modo di evitare 91
login per usare 1257	modo di fare 1396
luogo da visitare 64	modo di frequentare 65
macchina da scrivere 104	modo di giocare 117
macchina per scrivere 88	modo di incontrare 123
minaccia di fare 68	modo di intendere 206
minaccia di uccidere 95	modo di interpretare 78
modo da assicurare 87	modo di lavorare 173
modo da aumentare 84	modo di leggere 105
modo da avere 258	modo di mettere 113
modo da consentire 180	modo di operare 98
modo da creare 172	modo di ottenere 64
modo da dare 118	modo di parlare 300
modo da essere 126	modo di pensare 491
modo da evitare 248	modo di porre 108
modo da fare 386	modo di presentare 83
modo da favorire 65	modo di procedere 102
modo da formare 165	modo di produrre 83
modo da fornire 68	modo di raccontare 69
modo da garantire 187	modo di rendere 72
modo da impedire 77	modo di scrivere 187
modo da mantenere 104	modo di sentire 64
modo da migliorare 69	modo di studiare 84
modo da ottenere 300	modo di suonare 181
modo da permettere 253	modo di vedere 676
modo da potere 195	modo di vestire 153
modo da rendere 397	modo di vivere 590
modo da ridurre 148	modo per dire 74

-
- modo per evitare 66
modo per fare 328
modo per ottenere 85
modo per salvare 75
momento di fare 134
necessità di avere 223
necessità di costruire 104
necessità di creare 110
necessità di dare 90
necessità di fare 196
necessità di mantenere 87
necessità di trovare 110
necessità di utilizzare 75
obiettivo da raggiungere 65
obiettivo di creare 98
obiettivo di fare 128
occasione di conoscere 123
occasione di fare 119
occasione di parlare 75
occasione di vedere 83
occasione per fare 224
occasione per mettere 66
onore di essere 69
opportunità di fare 113
ora di vedere 87
ordine di uccidere 90
particolarità di avere 74
particolarità di essere 86
patto di avere 964
paura di essere 75
paura di fare 68
paura di perdere 209
permesso di costruire 91
possibilità di accedere 162
possibilità di acquistare 135
possibilità di andare 92
possibilità di aprire 76
possibilità di avere 390
possibilità di cambiare 120
possibilità di conoscere 121
possibilità di continuare 93
possibilità di controllare 82
possibilità di costruire 134
possibilità di creare 314
possibilità di dare 128
possibilità di diventare 138
possibilità di effettuare 236
possibilità di entrare 123
possibilità di eseguire 87
possibilità di esprimere 127
possibilità di essere 174
possibilità di fare 773
possibilità di giocare 260
possibilità di incontrare 67
possibilità di inserire 74
possibilità di lavorare 109
possibilità di mettere 132
possibilità di modificare 89
possibilità di montare 77
possibilità di operare 76
possibilità di ottenere 237
possibilità di partecipare 205
possibilità di passare 78
possibilità di portare 124
possibilità di prendere 93
possibilità di produrre 93
possibilità di raggiungere 110
possibilità di realizzare 177
possibilità di registrare 71
possibilità di ricevere 77
possibilità di ripubblicare 231
possibilità di salvare 108
possibilità di scaricare 67
possibilità di scegliere 362
possibilità di scrivere 66
possibilità di seguire 74
possibilità di sfruttare 93
possibilità di studiare 89
possibilità di sviluppare 94
possibilità di tornare 98
possibilità di trasformare 67
possibilità di trovare 115
possibilità di usare 222
possibilità di usufruire 91
possibilità di utilizzare 356
possibilità di vedere 314

- possibilità di vincere 148
possibilità di vivere 119
potere di fare 96
pretesa di essere 72
prezzo a partire 104
prezzo da pagare 166
probabilità di trovare 72
problema da risolvere 96
procinto di partire 66
procinto di sposare 90
prova di essere 69
punto di fare 82
punto di share 107
ragion d' essere 102
ragione d' essere 73
ragione di esistere 83
rischio di fare 106
rischio di perdere 131
rischio di sviluppare 78
riuscita a capire 114
scelta di fare 83
scopo di aiutare 77
scopo di assicurare 83
scopo di aumentare 127
scopo di costruire 64
scopo di creare 236
scopo di dare 132
scopo di difendere 92
scopo di diffondere 69
scopo di eliminare 103
scopo di evitare 181
scopo di facilitare 72
scopo di fare 459
scopo di favorire 110
scopo di fornire 138
scopo di garantire 105
scopo di impedire 92
scopo di mantenere 152
scopo di mettere 103
scopo di migliorare 112
scopo di ottenere 234
scopo di permettere 80
scopo di portare 87
scopo di preservare 69
scopo di produrre 83
scopo di promuovere 229
scopo di proteggere 160
scopo di raccogliere 111
scopo di realizzare 69
scopo di rendere 317
scopo di ridurre 163
scopo di studiare 64
scopo di trovare 88
scopo di visualizzare 485
sensazione di essere 92
sogno di diventare 199
speranza di fare 66
speranza di ottenere 94
speranza di riuscire 64
speranza di trovare 133
speranza di vedere 76
strada da percorrere 66
tempo di fare 161
tempo per fare 75
tentativo di conquistare 139
tentativo di costruire 89
tentativo di creare 210
tentativo di dare 177
tentativo di difendere 69
tentativo di distruggere 65
tentativo di evitare 73
tentativo di fare 418
tentativo di fermare 135
tentativo di liberare 90
tentativo di mantenere 65
tentativo di mettere 100
tentativo di migliorare 69
tentativo di ottenere 109
tentativo di portare 104
tentativo di prendere 65
tentativo di raggiungere 149
tentativo di realizzare 66
tentativo di recuperare 143
tentativo di rendere 121
tentativo di ridurre 67
tentativo di riportare 92

tentativo di risolvere 78	voglia di andare 83
tentativo di salvare 320	voglia di cambiare 64
tentativo di superare 87	voglia di fare 464
tentativo di trovare 149	voglia di giocare 77
tentativo di uccidere 109	voglia di lavorare 88
vantaggio di avere 77	voglia di vivere 231
vantaggio di essere 140	volontà di fare 137

Di seguito viene inoltre riportato l'insieme dei dati prodotti dallo strumento computazionale per tutte le espressioni considerate. I nomi delle colonne sono precisati dalla seguente legenda:

ESPRESSIONE	Espressione considerata
FREQ	Frequenza di occorrenza (componenti lemmatizzati)
FREQ_INT1	Numero di occorrenze dell'espressione interrotta tra primo e secondo componente
FREQ_INT2	Numero di occorrenze dell'espressione interrotta tra secondo e terzo componente
MOD_SINTAG	Valore di I_{syn}
FREQ_SIN1	Numero di occorrenze delle espressioni con primo componente sostituito (valore non intero poiché normalizzato in base alla similarità semantica)
DISP1	Dispersione della sostituzione del primo componente
FREQ_SIN2	Numero di occorrenze delle espressioni con secondo componente pieno sostituito (valore non intero poiché normalizzato in base alla similarità semantica)
DISP2	Dispersione della sostituzione del secondo componente
MOD_PARAD	Valore di I_{sub}
FORMA_PREVAL	Forma con il più alto numero di occorrenze
FREQ_P	Numero di occorrenze della forma prevalente
FORMA_PREVAL_2	Forma con il secondo numero di più alto di occorrenze
FREQ_P_2	Numero di occorrenze della seconda forma prevalente
MOD_FLESS	Valore di I_{infl}
ERR	Se contrassegnato da X, espressione scartata per errori dovuti al corpus

219	13	0	0.056034428	4.7098811022	0.94549667	986.64089444	0.840457957	0.8190607324	modo di esprimere	126 modo di esprimere	69	0.4246575342
217	0	2	0.056034428	208.2213296706	0.335791369	14.1874663635	0.335791369	0.5061540293	linea di evitare	210 linea di evitare	7	0.0326575342
215	1	1	0.0092165899	10.3114966757	0.128951844	721.987733197	0.1008206545	0.7730388842	grado di affrontare	200 grado di affrontare	15	0.0697674419
210	41	0	0.183446135	563.6565371519	0	10.365671508	0	0.071261782	165 tentativi di creare	155 tentativi di creare	47	0.2142857143
209	7	1	0.0368663594	0	0	0	0	0.0475259336	prezzo da perdere	159 prezzo da perdere	47	0.2392344498
207	6	1	0.0327102804	10.0336247719	0.185856573	375.22133093258	0.1200377365	0.6504889617	grado di mantenere	172 grado di mantenere	35	0.1690821256
206	6	4	0.0462962963	395.5203740885	0.0671595676	102.2456343314	1.3071662633	0.7272889617	grado di fare	113 grado di fare	58	0.4514563071
206	11	0	0.0506912442	3.4776351237	0.4145988835	586.6483272473	0.075840232	0.7412469763	modo di intendere	167 modi di intendere	33	0.1983203883
205	2	1	0.0142305738	5.312027026799	0.2071400218	1721.229788367	0.0053652784	0.8941691044	grado di operare	204 grado di operare	1	0.0048780488
205	2	1	0.0142305738	5.3105092973	0	306.5735889003	0.1154732298	0.6026239214	grado di svolgere	201 grado di svolgere	1	0.0195121951
205	10	0	0.0888888889	5.4552816769	0.3604776685	157.8055668037	0.3627460539	0.4433293909	modo di concepire	191 modi di concepire	14	0.0092968629
205	13	4	0.0765765766	0.9023094573	0.8408431335	92.2828906748	0.1362389724	0.3125077975	possibilità di partecipare	201 possibilità di partecipare	4	0.0195121951
200	2	0	0.0196078431	5.0056652216	0.7951481826	338.3486514769	0.1351410738	0.6319160558	grado di capire	168 grado di capire	32	0.16
199	1	0	0.006	2.9884925672	0	0	0	0.0148002328	associazione a delinquere	175 associazioni a delinquere	24	0.1206030151
199	1	4	0.0245098039	4.2080029109	0.6597608978	72.2819759578	0.7172245621	0.2776500331	grado di mettere	175 grado di mettere	28	0.1008040201
199	8	2	0.0478468927	114.6686182788	0.2036013663	501.2882193414	0.2005966577	0.7558715288	idea di creare	196 idea di creare	3	0.0150753769
199	6	0	0.0292682927	90.4652350349	0.1347151461	4.889097426	0	0.3239734715	scopo di diventare	196 scopo di diventare	3	0.0150753769
198	1	5	0.0294117647	293.4995309948	0.1206671056	54.4341406802	0.107720601	0.6373185787	linea di ottenere	191 linea di ottenere	6	0.0353535354
196	2	1	0.0150753769	0	0	24.57457426193	0.387766869	0.5561350305	grado di trasformare	103 grado di trasformare	93	0.4744897959
196	1	17	0.0841121495	28.9834766184	0.0814727253	368.7316195719	0.5332861937	0.6698753303	necessità di fare	90 necessità di fare	85	0.5408163265
195	5	0	0.025	0.6534672797	0	168.952448399	0.2668181982	0.4651363702	grado di soddisfare	190 grado di soddisfare	5	0.0256410256
195	0	2	0.010152843	0.6434384076	0	27.6748510624	0	0.1248592008	modo da potere	195	0	
190	1	1	0.0104166667	0.2555151022	0	237.4716699017	0.1012077045	0.5557916138	grado di ospitare	195 grado di ospitare	5	0.0263157895
189	0	1	0	0.179181056024	0.1540524013	0.2799618056024	0.1540524013	0.5984017842	grado di resistere	187 grado di resistere	1	0.0105820106
188	14	10	0.1132075472	139.38741228518	0.2022811692	176.5848673866	0.949177727	0.6273320099	diritto di fare	85 diritto di fare	61	0.5478723404
188	10	2	0.032673158	2.207328								

141	3	1.0506968214	2278.1956129183	0.2233935007	0.9417919698	grado di realizzare	133	grado di realizzare	6	0.0567375807	
140	1	0.0275862069	3.15057323414	0.0077071995	0.463.1396056046	0.7680212574	grado di comunicare	137	grado di comunicare	3	0.0214285714
140	2	0.0070921986	0.3677071995	0.6945008992	541.0312472744	0.189643241	0.7951178803	grado di esprimere	37	0.2642857143	
140	6	0.1907514451	0.0375700362	2.8831249505	5.5797339838	0.0622715779	vantaggio di essere	119	vantaggi di essere	21	0.15
139	31	0.1823529412	0.1824946672	0.299.935645212	0.251.6639161	0.683456798	tentativo di conquistare	114	tentativi di conquistare	11	0.1796561151
138	1	0.0071194246	0.2268884715	0.2268884715	0.0	0.0034967309	autorizzazione a procedere	110	autorizzazioni a procedere	28	0.202898507
138	6	0.0612244898	30.6456604899	0.7795545503	35.2226881342	0.2967878541	possibilità di diventare	138	possibilità di diventare	0	0
138	3	0.0349650535	57.6896731222	0.2422843755	29.3096713725	0.64808246	scopo di fornire	135	scopo di fornire	3	0.0217391304
137	8	0.1383647799	603.9387473757	0.195505987	112.3968764708	1.2535854708	volontà di fare	46	volontà di fare	46	0.6204379562
136	2	0.0285714286	333.4312872734	0.0889879243	48.9797770507	0.3731859413	linea di rendere	96	linea di rendere	40	0.2941176741
136	0	0.0144927536	0	0	0	0.25411203974	grado di contrastare	14	grado di contrastare	14	0.102941176
135	0	0	0	0	0	0.418.005009221	associazione per delinquere	134	associazioni per delinquere	1	0.0074074074
135	1	0.0145985401	0	0	0	0	possibilità di acquistare	130	grado di sopportare	15	0.030737037
135	5	0.0594440569	2.735664002	3.0105155504	431.1186524164	0.1601268242	possibilità di acquistare	123	possibilità di acquistare	12	0.0888888889
135	16	0.1176470586	0	0	0	0.156.551100805	tentativo di fermare	26	tentativo di fermare	26	0.2962362903
134	0	0.035971223	0.4369897908	1.9169350451	51.6310251756	2.2601631262	momento di fare	85	momento di fare	25	0.965716418
134	7	0.0694444444	48.761746414	0.1077950451	434.7694018468	0.1516852297	possibilità di costruire	122	possibilità di costruire	12	0.0895523388
133	3	0.0362318841	0.0643701869	0.4218570027	0.9164189345	0.0311979061	linea di essere	128	linea di essere	33	0.0376593985
133	6	0.0633802817	22.9713977934	0.5005656232	249.1397448812	0.0636490243	intenzione di lasciare	97	intenzione di lasciare	33	0.2706766937
133	9	0.0698300669	124.3095022895	0.1657843816	55.3577465505	1.3089086508	speranza di trovare	115	speranza di trovare	11	0.1353383459
132	1	0.007518797	0.0118114101	0	544.9750257582	0.4470710651	grado di spingere	111	grado di spingere	21	0.1509090091
132	1	0.0222222222	46.1471093597	0.2925971072	110.4545770324	0.3715855324	possibilità di mettere	87	possibilità di mettere	44	0.3409090909
132	4	0.0434762609	135.5566332276	0.1521781963	17.5544856196	1.32398642	scopo di dare	118	scopo di dare	40	0.1060606061
131	1	0.0507246377	110.1165896178	0.12923638196	12.923638196	0.1566749562	rischio di perdere	96	rischio di perdere	24	0.2671755725
129	4	0.0373134328	121.4649660895	0.1817681022	59.6207236644	0.2094503502	linea di garantire	113	linea di garantire	20	0.0775193798
129	3	0.0373134328	0.6758466473	0.6758466473	0.685.237749626	0.0558341762	grado di reggere	123	grado di reggere	6	0.0465116279
128	1	0.0076923077	2.0023748095	0	259.3294777458	0.2250518964	grado di uccidere	121	grado di uccidere	8	0.0620155039
128	1	0.0393037009	0.3368677831	0.7047498256	696.341575013	0.3862345023	grado di ottenere	125	grado di ottenere	3	0.023

111	3	0	0.0263157695	32.8762228455	0.2959116376	1.5334716963	0.285442767	scopo di raccogliere	109	scopo di raccogliere	2	0.018018018
110	14	0	0.1290322598	4219.7310519737	0.090426688	184.2299437928	0.9273429302linea di far	60line di fare	60	linea di fare	28	0.4945454545
110	6	1	0.0592905988	129.4272898473	0	280.4967811502	0.7465373961necessità di creare	105necessità di creare	105	necessità di creare	25	0.0454545455
110	1	0	0.0090909091	8.1954911185	0	254.123389821	0.7639126028necessità di trovare	104possibilità di trovare	104	possibilità di trovare	2	0.0272727273
110	8	1	0.0756302521	27.3159393067	0.2953602012	328.6131793155	0.0680960808	0.7639126028possibilità di raggiungere	104	possibilità di raggiungere	6	0.0545454545
110	5	0	0.0434782609	53.1324229642	0.162906744	450.4116701081	0.0956188616	0.7639126028scopo di lavorare	104	scopo di lavorare	3	0.0272727273
110	4	0	0.0801801801	40.4025989804	1.7780563208	347.4940332042	0.1017671988	0.76143748716grado di modificare	104	grado di modificare	5	0.0458715596
109	2	2	0.0353982301	17.410394931	0	516.5508310671	0.3640370422	0.830418923possibilità di lavorare	107	possibilità di lavorare	2	0.0183486239
109	32	0	0.2695903546	0.1515217381	0	311.9852173084	0.9097693986	0.7409303468tentativo di ottenere	78	tentativo di ottenere	22	0.20181348624
109	12	0	0.0991735537	0.15057021	0	131.1953205327	0.4040608913	0.5462026498tentativo di uccidere	76	tentativo di uccidere	26	0.3577981651
108	3	1	0.0608695652	1.3030608633	0.1030608633	108.1700595165	0.1020473342	0.5462026498tentativo di uccidere	76	modi di porre	26	0.2962962963
108	7	1	0.0357142857	9.2605372455	0.8045709788	93.5810674436	0.4166773246	0.5462026498tentativo di uccidere	76	possibilità di salvare	44	0.4074074074
107	4	1	0.046428571	0.3507321	0	138.510674436	0.4166773246	0.5462026498tentativo di uccidere	61	punto di share	46	0.4299056421
107	1	0	0.092592593	1.216570304	0.7227637366	351.5881715244	0.170914263	0.7672925384grado di costruire	98	grado di costruire	4	0.25
107	1	1	0.0183486239	17.9346976469	0.1421872512	312.9814765585	0.1390359921	0.7556069946grado di vincere	104	grado di vincere	2	0.0280373832
106	1	1	0.0093457944	2.7710316652	0.7567862261	146.18228805	0.1038134432	0.7556069946grado di vincere	50	grado di autare	50	0.1769881132
106	1	1	0.0093457944	3.0404950747	1.510151528	163.9813530741	0.2110151345	0.6117526831grado di mostrare	96	grado di mostrare	14	0.1320755477
106	4	2	0.0535714286	5.0091127747	0	16.5146795738	2.63236484516	0.6117526831grado di mostrare	34	fischio di fare	33	0.679245283
105	4	0	0.0366972477	0	0	338.8354181605	0.137951324	0.7672425022modo di leggere	94	modo di leggere	8	0.1047619048
105	1	0	0.0094339623	139.0551488924	0.0466861832	246.2789168555	0.0803730077	0.7858602791scopo di garantire	95	scopo di garantire	10	0.0952380952
104	3	2	0.0458715596	469.4540984402	0.0758045563	194.5792339891	0.3663934858	0.8645892808tentativo di creare	99	intento di creare	5	0.04807689231
104	0	0	0.0370377037	0.15125072364	0	78.615883665	0.2822718727	0.8645892808tentativo di creare	71	macchina da scrivere	33	0.3173076923
104	1	1	0.0188679245	0.402190355279	0.1650278815501	0.2395282537	0.6325734749necessità di costruire	91	modo da mantenere	13	0.125	
104	0	0	0.1186440678	3.5657854549	0	0	0.03314885prezzi a partire	90	modo da mantenere	13	0.125	
104	24	0	0.1875	0.61230408569	0	161.2923222315	0.0460077393	0.03314885prezzi a partire	80	prezzo a partire	24	0.0384615385
104	14	0	0.028301886									

permissio di costui-	3	3	1	0421052632	3.1513841116	0.5299190018	10.3267523075	1.1940258303	0.1290043724	permissio di costui-	36	5	0549450549
possibilità di usufrui-	91	5	1	0400000008	1.4748687878	0	526.559268195	0.0238666668	0.8529.18932	possibilità di usufrui-	90	5	0100890111
dovere di fare	5	2	0.0712649485	237.815011165	0.069737156	10.3135496968	4.77210125522	0.699854513	0.732920508	dovere di fare	29	29	0511111111
line di ridurre	4	0.0526315789	200.1752172289	0.1680426809	0.1680426809	0.1680426809	0.1680426809	0.682066278	0.4481908131	line di ridurre	14	15	0155555556
grado di costui-	90	1	0.010969011	0.3814528115	1.0909726567	133.456176498	24.151388182	0.5979228324	0.9213388182	grado di costui-	75	75	0166666667
nessità di dare	2	5	0.0721649485	31.4167839928	1.5721804516	63.2894811603	3.1493300648	0.5127399633	0.9213388182	nessità di dare	71	0	0177777778
ordine di uccide-	90	4	0.0526315789	1.2788634132	1.5721804516	58.408792078	0.7180240228	0.3987480150	0.9213388182	ordine di uccide-	61	61	0122222222
prointo di sposare	90	0	1.010969011	0	0	2.9876878821	3.486984502	0.0321320116	0.9213388182	prointo di sposare	20	20	0222222222
lervativo di libere	90	11	0.1089108911	0	0	31.5042930414	2.690133591	0.2592854373	0.9213388182	lervativo di libere	33	33	0444444444
ditto di vivere	89	3	0.0531914894	2.40161442736	1.32092424729	47.7323681713	0.1876136605	0.3602869707	0.9213388182	ditto di vivere	57	57	0.0224719101
modo di affrontar-	89	26	1	0.2327580627	2.2161662346	1.8516560455	190.2155296699	0.1565441011	0.6837598561	modo di affrontar-	13	13	0235959506
possibilità di modifi-	89	2	1	0.032606957	32.8366317198	0.3446332437	221.5894457113	0.1952951288	0.7408467037	possibilità di modifi-	68	68	0.0112359506
possibilità di modifi-	89	2	1	0.032606957	32.8366317198	0.3446332437	221.5894457113	0.1952951288	0.7408467037	possibilità di modifi-	68	68	0.0112359506
possibilità di studiar-	89	2	0	0.021978022	3.2167800616	1.0510778075	182.3452766928	0.2923300944	0.6758479479	possibilità di studiar-	90	90	0.0112359595
lervativo di costruire	89	14	0	0.1359223301	0	0	236.678287895	0.1767855596	0.7266402481	lervativo di costruire	62	62	0.0303707865
desiderio di diventare	88	7	2	0.0927835052	214.61798952	0.0607415922	6.1248471799	0.4071145019	0.7149731442	desiderio di diventare	88	88	0
grado di valutar-	88	0	0	0	0	0	309.4704203527	0.1616578972	0.7790937069	grado di valutar-	83	83	0.0566818182
macchina per scrivere	88	0	0	0	0	0	10.727681995	0.1196062681	0.9213388182	macchina per scrivere	39	39	0.44318182
scopo di trovare	88	2	2	0.0434782609	18.0905263316	1.384530682	0.6846214855	0.5066818182	0.2614062103	scopo di trovare	83	83	0.0566818182
volgia di lavorar	88	1	3	0.0434782609	38.6170371565	0.2539593046	115.8536519829	0.161950467	0.6846214855	volgia di lavorar	84	84	0.0566818182
modo da assicurare	87	0	0	0	0.0159423344	0	0.2327580627	0.227065597	0.7285536297	modo da assicurare	62	62	0.0476104716
nessità di mantenere	87	1	1	0.0224719101	7.9076869001	0	120.6272619597	0.333664163	0.5963451116	nessità di mantenere	25	25	0.2873563218
ona di vedere	87	1	3	0.0439560044	88.654677218	0.1030014018	25.7609458152	0.871155346	0.5680611797	ona di vedere	35	35	0.402298506
possibilità di eseguir-	87	2	0	0.0224719101	77.60061789974	0.1529577988	112.8320494475	0.0476199535	0.9296432645	possibilità di eseguir-	68	68	0.0566818182

Capacità di dare	67	4	0	0.056330282	0	0	0	161.8956445622	1.2349954891	0.7072901928	Capacità di dare	56	capacità di dar	8	0.1641791045
grado di indurre	67	1	2	0.0428571429	0	0	0	240.7770967243	0.144647582	0.782309989	grado di indurre	65	grado di indur-	2	0.0298597463
grado di raccogliere	67	0	0	0	0	0	0	157.5474693035	0.2237317953	0.7016221104	grado di raccogliere	56	grado di raccogli-	1	0.0149253731
possibilità di incontrare	67	1	1	0.0289855072	3.6321736378	0.3811152077	0	166.1786766262	0.29792900198	0.7194432333	possibilità di incontrare	51	possibilità di incontrar-	16	0.2388059701
possibilità di scaricare	67	4	1	0.0694444444	1.6486964109	0.1404885746	0	30.449167373	0.6079310207	0.323900663	possibilità di scaricare	61	possibilità di scaricar-	6	0.089522388
possibilità di trasformare	67	1	0	0.0147058824	98.7652384975	0.1404885746	0	195.1188847723	0.1138212313	0.814344839	possibilità di trasformar-	39	possibilità di trasformar-	28	0.4179104478
tentativo di ridurre	67	9	0	0.1184210526	0	0	0	6.1646839076	1.2066526714	0.084257643	tentativo di ridurre	57	tentativi di ridurre	8	0.1492537313
capacità di comunicare	66	7	2	0.12	2.999926183	0	0	109.3231596676	0.8497854036	0.6299026363	capacità di comunicare	66		0	0
capacità di operare	66	6	0	0.0633333333	0.8909440296	0	0	243.8466494133	0.0267637772	0.7876021396	capacità di operare	66		0	0
grado di condurre	66	2	2	0.0571428571	0.4356886542	1.9163484316	0	1076.5746385826	0.0921193268	0.942257278	grado di condurre	58	grado di condur-	8	0.1212121212
grado di eliminare	66	0	3	0.0434782609	0.888049943	1.21041815	0	272.0114250675	0.1586419706	0.8052519851	grado di eliminare	59	grado di eliminar-	7	0.1060606061
grado di visualizzare	66	1	0	0.049253731	0	0	0	183.8166364868	0.0399500842	0.735806226	grado di visualizzare	64	grado di visualizzar-	2	0.0303030303
modo di dimostrare	66	3	0	0.0434782609	0	0	0	431.3846478244	0.3026215953	0.8673059165	modo di dimostrare	50	modo di dimostar-	12	0.2424242424
modo per evitare	66	19	1	0.2325581396	20.0839258985	0.4625058568	0	15.4391927635	0.4076682934	0.3498657957	modo per evitar-	54	modo per evitar-	8	0.1818181818
occasione per mettere	66	7	0	0.095890411	8.5377315602	0.738067927	0	5.6593662857	2.1142409977	0.170275987	occasione per metter-	35	occasione per metter-	20	0.4698969697
possibilità di scrivere	66	0	1	0.0149253731	16.4827658731	0.1042846836	0	425.2475871461	0.9373878017	0.8700097398	possibilità di scrivere	65	possibilità di schver-	1	0.0151515152
proscinto di partire	66	0	0	0	0	0	0	0.6951346353	0.8558373224	0.0104225689	proscinto di partire	66		0	0
speranza di fare	66	1	5	0.0833333333	124.28222615731	0.1629416096	0	16.3950487872	2.7752303574	0.6806616078	speranza di far-	24	speranza di far-	24	0.6363636364
strada da percorrere	66	6	1	0.095890411	18.9244093602	0.2955112081	0	462.1984040401	0.0559926827	0.2228382806	strada da percorrere	52	strade da percorrere	14	0.2121212121
tentativo di realizzare	66	9	0	0.12	0	0	0	70.5433176565	1.588957892	0.8750469521	tentativo di realizzare	42	tentativi di realizzare	19	0.3636363636
bisogno di trovare	65	3	1	0.0579710145	120.2139500496	0.1121075911	0	38.870135714	1.6106743817	0.7458527745	bisogno di trovare	62	bisogno di trovar-	2	0.0461538462
desiderio di conoscere	65	2	0	0.0298507463	26.2929451777	0.4734882743	0	103.4760810285	0.9107298258	0.5006264483	desiderio di conoscere	53	desiderio di conoscer-	12	0.1846153846
diritto di avere	65	2	0	0.0298507463	15.2842539803	1.0541008643	0	32.9227190266	0.1851182584	0.6462783985	diritto di avere	62	diritto di aver-	2	0.0461538462
grado di passare	65	0	1	0.0298507463	2.3894901237	0.4108776906	0	62.2474678753	0.1531219503	0.4986384446	grado di passare	65		0	0
grado di prevedere	65	0	1	0.0151515152	0.2814222951	0	0	43.4013329146	0.3298462857	0.4003763768	grado di prevedere	61	grado di preveder-	3	0.0615384615
libertà di scegliere	65	0	0	0	0	0	0	29.9437369251	0.0889640085	0.8443773714	libertà di scegliere	64	libertà di sceglir-	1	0.0153846154
modo da favorire	65	3	2	0.0714285714	322.733310298	0.0223683856	0	108.4790822944	0.3983532345	0.6269518872	modo da favorir-	56	modo da favorir-	9	0.1384615385
modo di frequentare	65	0	3	0.0441176471	0.7611772808	0	0	305.1807550637	0.1442833636	0.8244101038	modo di frequentar-	63	modo di frequentar-	2	0.0307692308
obiettivo da raggiungere	65	4	1	0.0714285714	0	0	0	1.600667814	0.6238516479	0.127001098	obiettivo da raggiungere	42	obiettivi da raggiungere	23	0.3538461538
tentativo di distruggere	65	9	0	0.1216216216	7.8553260605	0.59538243	0	15.8063758008	1.2691835611	0.1956080278	tentativo di distruggere	54	tentativo di distrugger-	5	0.1692307692
tentativo di mantenere	65	12	0	0.1558441558	0	0	0	93.8511351022	0.1542785923	0.5908118632	tentativo di mantenere	53	tentativi di mantenere	7	0.1846153846
tentativo di prendere	65	8	0	0.1095880411	0	0	0	32.3615036479	0.6786481932	0.3323850026	tentativo di prendere	47	tentativo di prender-	8	0.2769230769
grado di scrivere	64	0	0	0.2073170732	0	0	0	978.1501291378	0.0350910991	0.9388334643	grado di scrivere	61	grado di scriver-	3	0.046875
intenzione di abbandonare	64	1	4	0.0724637681	4.173637161	0.3974601861	0	156.2483301054	0.1159710872	0.7276343517	intenzione di abbandonar-	57	intenzione di abbandonar-	7	0.109375
luogo da visitare	64	8	0	0.1111111111	14.7299074489	0.6997919526	0	0.5765510101	0	0.11998930462	luogo da visitare	59	luogo da visitare	5	0.078125
modo di ottenere	64	5	0	0.0724637681	8.150147025	1.179954512	0	148.11795950315	0.3477632661	0.6997038634	modo di ottenere	45	modi di ottenere	16	0.296875
modo di sentire	64	3	5	0.1111111111	0.4978582872	1.8829227761	0	141.6655997937	0.4437612073	0.6895667045	modo di sentir-	40	modo di sentir-	15	0.375
scopo di costruire	64	1	4	0.0724637681	37.8745749383	0.390268333	0	232.3975712444	0.1563206671	0.8085392375	scopo di costruir-	61	scopo di costruir-	3	0.046875
scopo di studiare	64	3	1	0.0588235294	14.0026251875	0.5494830797	0	72.1225266215	0.5588450714	0.5736890239	scopo di studiar-	58	scopo di studiar-	6	0.09375
speranza di riuscire	64	0	2	0.0303030303	6.0752239402	1.1721800872	0	71.1940678765	0.0880267614	0.5469645301	speranza di riuscire	58	speranze di riuscire	4	0.09375
voglia di cambiare	64	2	1	0.0447761194	32.500286847	0.219533773	0	8.7888917692	0.8847994584	0.3920347737	voglia di cambiar-	63	voglia di cambiar-	1	0.015625
bisogno di dire	63	1	2	0.0454545455	9.4297295593	1.586217543	0	112.8748699733	2.2224467849	0.6600194158	bisogno di dir-	34	bisogno di dir-	29	0.4603174603
grado di attivare	63	0	1	0	0	0	0	115.7113742202	0.1054468852	0.6474762713	grado di attivar-	56	grado di attivar-	7	0.1111111111
grado di camminare	63	0	1	0.015625	0.4563529898	0	0	246.5344504726	0.0915525505	0.7967681644	grado di camminare	63		0	0

G

Dati sul pattern NN

Si riporta di seguito la lista, in ordine alfabetico, delle 500 espressioni relative al pattern NN estratte in PAISÀ e i dati relativi al comportamento variazionale di ognuna delle espressioni.

- - 1049	ascolto tg 367
- - 806	ascolto tv 302
- br 273	asilo nido 282
- link 1191	assetto variabile 201
's anatomy 196	associazione ambientalista 294
a est 270	attacco aereo 414
acciaio inox 185	attacco nemico 249
acqua marina 187	attività live 216
acqua minerale 247	attore protagonista 1270
aereo militari 176	attore statunitense 214
aereo nemico 176	azienda statunitense 307
aiuto regista 235	bambino prodigio 213
albero motore 466	banca dato 348
album live 546	band punk 177
album solista 803	been deleted 244
album studio 144	been downloaded 244
already sent 190	before the 244
alter ego 705	best seller 323
an der 213	black bloc 128
and sisters 213	black metal 441
and the 205	blog personale 250
anno luce 1137	bomba atomica 174
anticipo rispetto 247	bombardamento aereo 248
anziano signore 199	bonus track 419
apparato motore 299	browser closed 244
applicazione web 194	busta paga 223
arco temporale 286	cache limiter 190
ascolto record 221	calciatore campione 1292

cambio nome 257
campagna acquisto 207
campionato costruttore 195
campo base 230
campo medico 386
campo profugo 362
candidato sindaco 246
cardinale vescovo 216
casa automobilistica 989
casa madre 588
casa natale 384
cassa integrazione 356
centro abitato 380
centro benessere 361
centro città 495
centro destra 332
centro sinistra 473
centro studio 203
cesare pavese 255
chiesa madre 405
chirurgia plastica 228
cinema horror 202
città greca 416
città lagunare 209
città natale 2600
classe dirigente 189
classe lavoratore 292
classifica costruttore 217
classifica marcatore 519
classifica pilota 179
closed the 244
codice sorgente 698
cofano motore 316
colonia greca 286
colore marrone 284
combattimento corpo 241
comitato organizzatore 173
commento a?? 385
commento giovedì 245
commento lunedì 206
commento martedì 236
commento mercoledì 197
commento venerdì 221

composizione chimica 448
computer grafica 294
comunicato stampa 732
condizione meteo 294
conferenza stampa 2370
connection before 244
cornice marcapiano 187
corpo vettura 545
could not 251
croce greca 501
cronologia presenza 583
cultura greca 295
dato ascolto 670
dato auditel 873
de la 450
de las 212
de los 471
decimo posto 265
decreto legge 437
dedicated to 249
deleted from 244
denied for 882
diario online 893
differenza rete 453
differenza rispetto 506
diretta conseguenza 209
diretta tv 217
direzione sud 192
domenica mattina 422
domenica pomeriggio 395
domenica sera 471
don 't 233
downloading this 244
dulcis in 134
effetto serra 479
eius et 189
elemento chiave 276
elenco agriturismo 204
elenco campeggio 203
elezione parlamentare 422
emergenza rifiuto 201
emisfero nord 2569
emisfero sud 2945

-
- episodio numero 423
 - episodio pilota 528
 - erede maschio 185
 - error message 2089
 - esercito nemico 222
 - esibizione live 341
 - esordio boom 2112
 - età avanzata 321
 - fabbricato viaggiatore 511
 - falso nome 315
 - fan club 191
 - fibra ottica 422
 - figlio femmina 242
 - figlio maschio 599
 - figlio minore 172
 - figlio primogenito 517
 - figura chiave 176
 - file sharing 232
 - file system 703
 - film commedia 2046
 - film documentario 288
 - film horror 1347
 - film tv 396
 - filo conduttore 346
 - filosofia greca 203
 - fine agosto 501
 - fine anno 1569
 - fine settembre 201
 - fine settimana 1279
 - fine stagione 724
 - floppy disk 300
 - flusso migratore 207
 - fondo pensione 252
 - for this 246
 - forma verbale 211
 - forum rete 541
 - forza lavoro 903
 - francesca camerino 341
 - fratello minore 814
 - from the 263
 - fuoco nemico 244
 - gaia scienza 252
 - galleria immagine 486
 - gas serra 468
 - general manager 175
 - geometria variabile 281
 - giornale online 367
 - giovane età 483
 - giovedì sera 234
 - giro @girodivite.it 231
 - grado centigrado 204
 - greatest hits 179
 - gruppo punk 371
 - guest star 688
 - hard disk 854
 - hardcore punk 276
 - has now 244
 - headers already 190
 - heavy metal 942
 - high school 219
 - home computer 234
 - home page 187
 - home video 536
 - home-mylifetv-public_html-classes-config.php on 380
 - imc-sf-active-shared-classes-db_class.inc on 4410
 - immaginario collettivo 534
 - in fundo 205
 - industria automobilistica 239
 - industria chimica 293
 - industria farmaceutica 225
 - information overload 292
 - informazione riguardo 254
 - inizio anno 452
 - inizio carriera 173
 - inizio secolo 261
 - inizio stagione 317
 - interfaccia utente 282
 - interrupted by 244
 - isola greca 214
 - it 's 249
 - it was 224
 - italiano - 362
 - la prima 196
 - largo anticipo 187

lato est 396
lato nord 629
lato ovest 364
lato sud 506
lavoratore dipendente 204
lettera greca 235
letteratura greca 204
lettore dvd 191
linea guida 6542
linea temporale 265
lingua greca 722
lingua madre 385
linguaggio macchina 198
link - 687
live action 197
live album 183
long playing 208
longer than 684
luna park 201
lunedì sera 239
lunghezza variabile 220
magister militum 310
maglia numero 410
mailing list 193
man mano 1912
martedì sera 200
mass media 412
materia prima 466
medical drama 417
membro fondatore 776
message above 244
message has 244
metà anno 195
metà classifica 231
metà stagione 351
metà strada 1121
metro cubo 1088
metro piano 329
metro s.l.m 209
mezzo secolo 190
migliori anno 277
miniserie tv 283
minor numero 390

mitologia greca 2218
modello base 207
modus operandi 199
mondo materiale 186
motore diesel 1079
motore elettore 330
movimento indipendentista 183
mura cittadino 377
nave ammiraglia 369
nord est 293
nord ovest 328
not send 244
notizia estero 246
now been 244
numero complesso 426
numero minimo 208
numero variabile 355
of the 452
offerta lavoro 526
oggi arcidiocesi 774
oggi diocesi 560
onda radio 480
open source 1220
opera lirica 603
origine greca 392
other possible 244
padre fondatore 318
paese membro 408
paese natale 829
pagina web 672
paolo polo 246
par condicio 176
parco divertimento 271
parco gioco 272
parete nord 286
parete sud 189
parola chiave 652
parola fine 203
parola greca 472
parroco don 220
parte est 197
parte nord 545
parte ovest 173

parte sud 356
particella carica 201
particolare tipo 704
partita casalinga 655
partito conservatore 237
partito socialista 326
pausa pranzo 180
pay tv 190
paziente affetto 303
performance live 235
personaggio chiave 191
pian terreno 237
pianeta natale 223
piano regolatore 350
piano terra 1666
piano terreno 380
pianta medicinale 267
pietra arenaria 239
pilota collaudatore 235
pit stop 319
polca position 930
polo nord 332
polo sud 244
popolo indigeno 187
portiere titolare 199
posizione numero 1013
potere temporale 562
prima nozze 222
principe elettore 342
programma tv 422
punto cardinale 321
punto chiave 247
punto critico 174
punto vendita 686
raccolta fondo 236
raggio gamma 384
raggio laser 215
rassegna stampa 197
razzo vettore 260
reality show 1305
reason for 244
reazione chimica 312
redattore capo 201
reddito pro 191
regina madre 236
regular season 669
reply was 922
reply within 244
request for 244
resource in 32460
result resource 8121
rete ammiraglia 263
ritardo rispetto 296
romanzo statunitense 537
ruolo chiave 330
russo - 240
sabato mattina 304
sabato pomeriggio 371
sabato sera 1522
sala gioco 219
sale minerale 309
santo patrono 657
scalo merce 379
scheda madre 560
scheda video 526
scienziato pazzo 240
secondo nozze 724
sede titolare 196
segnalazione inviato 224
segnale radio 202
segno opposto 203
send session 190
senso opposto 254
senso orario 568
serie anime 361
serie commedia 187
serie manga 304
serie tv 1768
server could 244
servizio passeggero 363
servizio viaggiatore 257
sesta stagione 361
sicurezza informatica 243
signor presidente 180
sistema binario 242
sistema radar 213

sito internet 1410
sito web 2812
smart card 193
soap opera 1131
socio fondatore 384
software house 327
sol colpo 176
soluzione alternativa 192
sostanza stupefacente 529
spada laser 469
specie animale 504
squadra avversario 259
squadra campione 522
squadra partecipante 1380
squadra riserva 184
stato membro 527
statunitense - 1163
statunitense vincitore 224
stazione radio 753
stella bianco-azzurre 222
stella binaria 196
stile neoclassico 697
struttura dato 286
successo grazie 248
sud est 343
sud ovest 340
tag team 371
talent scout 205
talent show 1420
talk show 955
team manager 238
temperatura ambiente 661
temperatura massima 227
tenente colonnello 669
terra natale 249
territorio nemico 238
test match 394
that the 258
the browser 244
the cache 244
the connection 244
the error 244
the request 246
the socket 244
the whole 244
this casa 683
this error 244
this url 244
thrash metal 276
time longer 244
title track 201
titolo costruttore 206
titolo pilota 189
to the 282
top ten 539
traccia audio 198
traccia bonus 292
traffico merce 316
traffico passeggero 405
trasmissione radio 192
trasporto merce 266
trasporto passeggero 242
trasporto truppa 343
trattamento medico 179
treno merce 349
truppa russo 175
tubo lanciasiluri 199
tv via 191
uccello migratore 203
ufficio stampa 387
using password 8975
valid mysql-link 2646
van der 558
vano scala 175
vantaggio rispetto 246
venerdì santo 190
venerdì sera 406
versante nord 232
versante sud 202
versione alternativa 371
versione base 407
versione berlina 277
versione live 338
versione online 356
vescovo titolare 201
via cavo 177

via de 303	whole of 244
via internet 115	within the 244
via mail 179	you can 156
via via 239	zona nord 295
video tributo 204	zona punto 264
vita cittadino 203	zona retrocessione 213
volta campione 453	zona sud 262
was interrupted 244	
week end 288	

Di seguito viene inoltre riportato l'insieme dei dati prodotti dallo strumento computazionale per tutte le espressioni considerate. I nomi delle colonne sono precisati dalla seguente legenda:

ESPRESSIONE	Espressione considerata
FREQ	Frequenza di occorrenza (componenti lemmatizzati)
FREQ_INT	Numero di occorrenze dell'espressione interrotta
MOD_SINTAG	Valore di I_{syn}
FREQ_SIN1	Numero di occorrenze delle espressioni con primo componente sostituito (valore non intero poiché normalizzato in base alla similarità semantica)
DISP1	Dispersione della sostituzione del primo componente
FREQ_SIN2	Numero di occorrenze delle espressioni con secondo componente sostituito (valore non intero poiché normalizzato in base alla similarità semantica)
DISP2	Dispersione della sostituzione del secondo componente
MOD_PARAD	Valore di I_{sub}
FORMA_PREVAL	Forma con il più alto numero di occorrenze
FREQ_P	Numero di occorrenze della forma prevalente
FORMA_PREVAL_2	Forma con il secondo numero di più alto di occorrenze
FREQ_P_2	Numero di occorrenze della seconda forma prevalente
MOD_FLESS	Valore di I_{infl}
ERR	Se contrassegnato da X, espressione scartata per errori dovuti al corpus

144	album studio	5	0.033557047	0	0.2316401239	1.5994689541	0.00106007	album studio	144	0	0
588	case madre	39	0.0022090569	0.4126424245	0	0.4026611954	0	0.0013846509	case madre	10	0.02040061033
583	cronologia presenza	0	0	0	0	0	0	0	583	0	0
589	senso orario	7	0.012173913	0.8999611914	0	0	0	0.0015819322	senso orario	1	0.00176656342
562	potere temporale	12	0.0209059233	11.8995852284	0	0.2024874793	0	0.0207175483	potere temporale	26	0.048263452
560	oggi diocesi	0	0	0	0	0	0	0	560	0	0
560	scheda madre	0	0	0	0	0	0	0	371	0.3375	0
558	van der	0	0	0	0	0	0	0	558	0	0
545	parte nord	2	0.00365563071	2.4086792038	0	0	0	0.004400148	parte nord	2	0.00366697248
545	corpo vettura	0	0	0.1066241287	0.4927758296	109.4561345233	0	0.1673831229	corpo vettura	19	0.0348623853
541	forum rete	0	0	0	0	0	0	0	541	0	0
539	top ten	0	0	0	0	0	0	0	539	0	0
537	romanzo statunitense	2	0.0037105751	0	0	0	0	0	537	0	0
534	immaginario collettivo	0	0.0092764378	0	0	0	0	0	534	0	0
529	sovrastanza stupefacenti	5	0.0093632969	0	0.0311789056	0	0.4319682256	0.0008159038	sovrastanza stupefacenti	3	0.0056710775
528	episodio pilota	1	0.0018903592	80.6245207195	0	0	0	0.13247005	episodio pilota	22	0.041666667
527	stato membro	12	0.0222634508	20.3574803028	0.1837346356	0	0	0.0371923163	stato membro	259	0.4990512384
526	offerta lavoro	1	0.0018975322	0	0	0	0	0	526	0	0
526	scheda video	6	0.0012781855	0	0	0	0	0	338	0.36121673	0
522	Squadra campione	22	0.003887716	0.1697906749	0	0	0.24255625	0.1882806546	Squadra campione	43	0.0842918177
519	classifica marcatori	2	0.003887716	0.1697906749	0	0	0	0.1882806546	classifica marcatori	510	0.0173410405
517	figlio primogenito	29	0.0531135531	0	0	0	0	0	517	0	0
511	differenza rispetto	506	0	7.6414851832	0	4.2532265165	0	0.0227478141	differenza rispetto	239	0.4723320158
506	lato sud	151	0.2298325723	53.7736926802	1.7040223562	8.1751581181	0	0.1090747005	lato sud	22	0.0434782009
504	specie animale	504	6	0.01171875	884.449244256	0.0574679027	0	0.6360888381	specie animale	484	0.0158730159
504	fine agosto	24	0.0454545455	26.4436635578	0.7232088633	0	0	0.0498519737	fine agosto	10	0.0199600798
501	croce greca	1	0.0019920319	0	0	0	0	0	501	0	0
501	croce greca	5	0.0098814229	0	0	0	0	0	491	0	0
495	centro città	32	0.0607210626	0	0	346.6891127733	0	0.411896872	centro città	10	0.0202020202
486	galleria immagine	1	0.0020533891	0	0	0.894427191	0	0	486	0	0
483	giovane età	4	0.0082135524	0	0	0.5337586446	0	0.0011038705	giovane età	483	0
480	onda radio	1	0.0020790021	0	0	0	0	0	448	0	0
479	effetto serra	2	0.0041580042	0	0	0	0	0	479	0	0
473	centro sinistra	26	0.0521042084	0	0	0	0	0	473	0	0
472	parola greca	10	0.020746898	0	0	0	0	0	316	0.3305084746	0
471	de los	2	0.0042780239	0	0	0	0	0	471	0	0
471	domenica sera	3	0.0063291139	0	0	0	0	0			

241	8	0.0321.2851.41	87.0980986395	0.2953390912	0	0	0.2654635895	combattimento corpo	203	combattimenti corpo	38	0.15.67/63495
235	13	0.0524.193549	326.279302189	0.0154650268	0	0	0.5913128726	performance live	235	0	0	0
240	0	0	0	0	0	0	0	russo -	240	0	0	0
240	6	0.0243902439	0	0	0	0	0	scienziato pazzo	173	scienziati pazzi	67	0.279166667
238	1	0.004184.1004	0	0	0	0	0	clean manager	238	0	0	0
239	7	0.02949552846	171.00447631142	0.0949885687	0	0	0.4171368068	industria automobilistica	239	0	0	0
239	0	0	0	0	0	0	0	pietra arenaria	239	0	0	0
239	12	0.0478087649	24.32061912221	0	0	0	0.0923612408	pietra arenaria	239	pietre arenarie	2	0.0083682008
239	1150	0.8279337653	0.6559704028	4.1757977979	0	0	0.0027367175	via via	239	0	0	0
238	3	0.01246301329	0.1068101519	10.812348321	0.3846829767	0	0.1243197009	territorio nemico	224	territori nemici	14	0.0588235294
237	51	0.1770833333	58.9492328232	0.3228381494	0.4916206927	0	0.2005150532	partito conservatore	198	partiti conservatori	39	0.164556962
237	0	0	0	0	0	0	0	pian terreno	237	0	0	0
236	0	0	0	0	0	0	0	commenti martedì	222	commento martedì	14	0.0593220339
236	2	0.0084033613	0	0	0	0	0	raccolta fondi	232	raccolte fondi	4	0.0169491525
236	14	0.0396	0	0	0	0	0	oregina madre	232	regine madri	4	0.0169491525
235	0	0	0	0	0	0	0	auto regista	235	0	0	0
235	11	0.0447154472	31.3428114704	0	0	0	1.0358114658	auto regista	107	pietre greche	107	0.4595744681
235	18	0.0711462451	0	0	0	0	0	lettera greca	174	pietre greche	59	0.2595744681
234	0	0	0	0	0	0	0	pilota collaudatore	234	0	0	0
234	0	0	0	0	0	0	0	ipocrita sera	174	pioti collaudatori	39	0.2010309278
194	4	0.0292020202	0	0	0	0	0	applicazioni web	155	applicazione web	39	0.2010309278
232	2	0.0085106393	0	0	0	0	0	don 1	233	0	0	0
232	0	0	1235.8067074386	0.014770502	0	0	0.8420486516	versante nord	220	versanti nord	12	0.0517241379
232	0	0	0	0	0	0	0	olle shairig	232	0	0	0
231	0	0	0	0	0	0	0	olgio @grodvile.it	231	0	0	0
231	0	0	50.6341849199	0	0.3227606584	0	0.180725981	meta classifica	231	0	0	0
230	2	0.0086206897	8.65696632	3.6053542692	0.1089512857	1.4950132784	0.0367130341	campo base	221	campi base	9	0.0391304348
224	1	0.0044444444	0	0	0	0	0	olt was	224	0	0	0
228	0	0	0	0	0	0	0	olitura plastica	228	0	0	0
227	16	0.0658436214	0	0	0	0	0	oltemperature maxsime	226	temperatura maxsime	1	0.0044052863
225	11	0.0466101695	43.2506593432	0.2386625478	0	0	0.1612322574	industria farmaceutica	224	industrie farmaceutiche	1	0.0044444444
224	0	0	1.6488846689	0.9406700308	0	0	0.0073073193	segnalazioni inviate	224	0	0	0
224	2	0.0088495575	0	0	0	0	0	ostatutensivi vincitori	224	0	0	0
223	1	0.0044642857	0	0	0	0	0	obusia paga	166	busste paga	57	0.2556053812
223	0	0	290.6334443995	0.0150605169	0	0	0.5658382404	planetaria natale	223	0	0	0
222												

comitato organizzatore	173	0	0	0	0	26.1307165987	0.1605443251	0.1312239369	comitato organizzatore	159	comitati organizzatori	10	0.0809248555
inizio carriera	173	17	0.0894736842	0	0	0.7952465934	1.1683190155	0.0045757672	inizio carriera	171	inizi carriera	2	0.0115606936
figlio minore	172	3	0.0171428571	0	0	6.0311470838	0.33419965	0.0338769209	figlio minore	171	figli minori	1	0.0058139535

H

Dati sul pattern NCN

Si riporta di seguito la lista, in ordine alfabetico, delle 500 espressioni relative al pattern NCN estratte in PAISÀ e i dati relativi al comportamento variazionale di ognuna delle espressioni.

acqua e sapone 71
acquedotto e impianto 67
adulto e bambino 79
aereo ed elicottero 80
agosto e dicembre 244
agosto e settembre 122
agricoltura e allevamento 65
alesaggio e corsa 206
alimento e bevanda 59
allucinazione e-o diceria 246
amico e collaboratore 140
amico e collega 143
amico e compagno 278
amico e conoscente 121
amico e familiare 69
amico e nemico 79
amico e parente 301
andata e ritorno 1087
anima e corpo 199
animale e pianta 77
animale e vegetale 93
anime e manga 378
anno e anno 274
anno e mezzo 1477
anno ed anno 56
apertura e chiusura 161
approfondimento e lettura 73
aprile e giugno 54

aprile e maggio 182
architetto e ingegnere 52
architettura e patrimonio 88
arco e freccia 91
arma e armatura 64
arma e munizione 150
arrivo e partenza 66
arte e cultura 76
arte e mestiere 175
arte e scienza 54
articolo e saggio 71
artista e gruppo 65
artista e intellettuale 61
artista e letterato 75
ascolto e spazio 61
aspirazione e scarico 54
asta e bilanciere 98
attività e diffusione 525
attore e cantante 141
attore e doppiatore 53
attore e regista 227
audio e video 179
auto e moto 70
autore e conduttore 120
autore e produttore 90
autore e regista 118
bambino e adolescente 53
bambino e adulto 61

- bambino e ragazzo 210
bar e ristorante 68
bene e servizio 386
bene o servizio 124
bianco e nero 61
biografia e carriera 148
biografia e opera 61
biologo e scrittore 93
botta e risposta 148
braccia e gamba 91
buono e cattivo 69
caccia e pesca 57
caccia e raccolta 54
calcio e magnesio 78
calcio e pugno 113
calice e corolla 84
campionato e coppa 90
cane e gatto 164
cantante e attore 85
cantante e chitarrista 196
cantante e compositore 66
cantante e gruppo 58
cantante e musicista 82
canto e ballo 96
canto e pianoforte 60
capo e promontorio 70
cappa e spada 76
caratteristica e informazione 66
carbonio e ossigeno 56
carico e scarico 125
carne e ossa 155
carne ed ossa 260
carta e penna 116
cattolico e protestante 52
causa ed effetto 117
centinaio e centinaio 125
chiesa e convento 183
chiesa e monastero 141
chiesa e palazzo 55
chilometro e mezzo 119
chitarra e voce 161
chitarrista e cantante 96
cibo e acqua 91
cibo e bevanda 90
cielo e terra 73
cinema e teatro 94
cinema e televisione 166
cinema e tv 128
citazione e riferimento 75
città e paese 60
città e villaggio 86
collega e amico 100
collega ed amico 59
colore e simbolo 118
comando e controllo 88
compagno e compagno 108
compositore e direttore 77
concetto e principe 173
conoscenza e accettazione 481
consumo ed emissione 116
coro e orchestra 78
cosa e persona 59
critica e pubblico 109
cultura e civiltà 132
cultura e religione 60
cultura e tradizione 62
dans ma tribu 111
decina e decina 379
decina o centinaio 62
decollo e atterraggio 79
decollo ed atterraggio 57
deputato e senatore 99
descrizione e stile 484
destra e sinistra 294
destra o sinistra 75
devastazione e saccheggio 101
dicembre e gennaio 54
dimensione e forma 54
dimensione e peso 98
diritto e dovere 226
diritto e libertà 54
domanda e risposta 111
domenica e lunedì 81
donna e bambino 490
donna e macchina 219
donna e uomo 248

-
- economia e commercio 113
economia e felicità 121
economia e manifestazione 162
edificio e monumento 64
emozione e sentimento 77
ente e associazione 60
ente ed istituzione 89
episodio e programmazione 111
est e ovest 81
est ed ovest 56
evento e manifestazione 107
falce e martello 128
familiare e amico 58
fatto e notizia 99
fauna e flora 85
febbraio e giugno 294
febbraio e marzo 135
fede e ragione 70
ferro e fuoco 335
ferro e vetro 55
festa e fiera 55
festa e sagra 72
figlio e figlio 66
figlio e nipote 130
figlio e successore 161
figlio ed erede 150
film e serie 76
film e telefilm 104
filosofia e teologia 133
fiume e lago 82
fiume e torrente 86
flora e fauna 405
fondatore e direttore 124
fondatore e presidente 102
fonte e bibliografia 116
fonte e collegamento 316
forma e colore 101
forma e contenuto 63
forma e dimensione 292
forma e-o contenuto 480
forza e resistenza 52
foto e video 87
fratello e sorella 391
fretta e furia 212
frutta e verdura 351
fumetto e cartone 100
gas e polverio 87
genere e numero 59
genitore e figlio 118
gennaio e febbraio 196
geografia e clima 78
giacca e cravatta 145
gioia e dolore 55
giornale e rivista 175
giornalista e scrittore 154
giorno e giorno 100
giorno e mezzo 77
giorno e notte 353
giugno e luglio 184
giugno e novembre 171
giugno e settembre 68
giustizia e libertà 79
greca e latino 161
greco e latino 80
gruppo e musicista 1785
gruppo e organizzazione 158
guelfo e ghibellino 251
guerra e pace 58
hardware e software 126
idrogeno ed elio 60
imbroglione e mezzo 79
immagine e somiglianza 152
immagine o video 119
incidente e disastro 76
infanzia e adolescenza 115
infanzia e gioventù 76
infanzia e giovinezza 94
infrastruttura e trasporto 869
inizio e fine 90
istituzione e carica 74
lago e fiume 61
latino e greco 108
lavoratore e lavoratore 60
legge e regolamento 71
legge o legge 67
lettera e filosofia 129

- lettura e scrittura 86
libertà e democrazia 57
libro e articolo 66
libro e film 71
libro e rivista 65
lingua e cultura 178
lingua e dialetto 248
lingua e letteratura 193
litro e mezzo 75
luce e ombra 96
luce ed ombra 122
luglio e agosto 322
luglio e massimo 93
luglio e settembre 53
luglio ed agosto 81
madre e figlio 268
maggio e giugno 159
maggioranza e opposizione 84
maglie e sponsor 105
maltrattamento o trattamento 73
mamma e papà 180
manga e anime 371
manga ed anime 99
manifestazione ed evento 72
mano e piede 135
marito e moglie 243
marzo e agosto 139
marzo e aprile 116
maschio e femmina 368
maschio o femmina 62
massa e potere 247
massima e minima 100
matematica e fisica 104
matrimonio e discendenza 223
matrimonio e figlio 362
matrimonio ed erede 144
medicina e chirurgia 149
medico e infermiere 59
medico e paziente 73
merce e passeggero 66
mese e mese 146
mese e mezzo 444
mese o anno 58
metafisica e quotidianità 246
metro e largo 135
metro e mezzo 432
migliaio e migliaio 182
miliardo e mezzo 143
milione e mezzo 1832
milione e milione 118
mini-serie o film 101
miniserie o film 146
minuto e mezzo 190
mito e leggenda 73
moglie e figlio 303
moglie e madre 73
moglie e marito 53
mondo e cultura 57
moneta e banconota 94
monumento e luogo 2631
morte e distruzione 136
morto e ferito 226
museo e galleria 56
musica e ballo 74
musica e danza 59
musica e potere 247
musica e testo 74
musicista e cantante 74
musicista e compositore 91
nascita e morte 57
nascita e sviluppo 54
nemico ed alleato 121
nome e città 481
nome e cognome 493
nord e sud 248
notizia e curiosità 108
notte e giorno 59
novembre e dicembre 170
operaio e contadino 61
ora e mezzo 176
ora e ora 168
ora ed ora 92
ordine e grado 126
organo e tessuto 52
origine e diffusione 147
origine e giacitura 166

-
- origine e storia 141
origine e sviluppo 55
origine ed evoluzione 52
oro e argento 219
oro ed argento 70
ottobre e aprile 253
ottobre e novembre 128
pace e prosperità 67
padre e figlio 472
padre e madre 63
pagina e pagina 66
parco e giardino 90
parente e amico 142
parente ed amico 63
parola e frase 68
parola e musica 55
partenza ed arrivo 70
partito e movimento 77
passeggero e merce 122
passo e valico 447
persona o cosa 52
personaggio e situazione 59
peso e misura 56
petrolio e gas 130
pianoforte e orchestra 179
pianta e animale 102
pianta e fiore 52
pianta ed animale 62
pilota e costruttore 82
pittore e scultore 124
pittura e scultura 142
poesia e politica 493
poeta e scrittore 187
polizia e carabiniere 84
poliziotto e carabiniere 84
porta e finestra 137
positività o critica 70
posizione e teoria 65
potenza e coppia 69
potere e abilità 683
pranzo e cena 55
pregio e difetto 127
premio e candidatura 67
premio e nomination 75
premio e riconoscimento 1146
presenza e rete 1285
presidente e amministratore 87
pressione e temperatura 109
primavera e inizio 99
primavera ed estate 76
prodotto e servizio 145
prodotto o servizio 95
produttore e regista 76
produzione e commercializzazione 62
produzione e distribuzione 178
produzione e vendita 53
profilo e caratteristica 100
profilo e storia 97
progettazione e costruzione 106
progettazione e realizzazione 93
promozione e retrocessione 71
protone e neutrone 66
pubblico e critica 289
punto e mezzo 53
punto e virgola 102
qualità e quantità 85
quantità e qualità 123
quartiere e frazione 65
quotidiano e periodico 62
quotidiano e rivista 102
racconto e romanzo 88
radice e fusto 66
radio e televisione 114
radio e tv 74
ragazzo e ragazzo 282
re e regina 74
realtà e finzione 72
regista e attore 81
regista e produttore 88
regista e sceneggiatore 138
religione e politica 137
ricerca e soccorso 62
ricerca e sperimentazione 54
ricerca e studio 64
ricerca e sviluppo 294
riconoscimento e premio 72

- riga e colonna 55
riga e compasso 112
rischio e pericolo 75
rivista e giornale 71
rivolta e rivoluzione 123
romanzo e racconto 211
sabato e domenica 257
saggio e articolo 58
sala e pepe 181
salita e discesa 54
sano e salvo 125
scena e costume 59
sceneggiatore e regista 84
scienza e coscienza 60
scienza e tecnologia 125
scrittore e critico 59
scrittore e giornalista 194
scrittore e poeta 188
scrittore e regista 59
scultore e architetto 52
scuola e corrente 90
scuola e università 62
secolo e mezzo 248
segno e sintomo 114
seno e coseno 60
serie e competizione 107
settembre e febbraio 315
settembre e ottobre 101
signore e signore 74
sintomo e segno 67
sito e-o blog 478
società ed istituto 234
soggetto e oggetto 84
soggetto e sceneggiatura 112
sogno e realtà 53
somma e prodotto 68
soprannome o nomiglioli 64
sorriso e canzone 71
spazio e tempo 142
squadra e corridore 59
stagione e mezzo 61
status e conservazione 139
stazione e comprensorio 254
stazione e fermata 152
stella e striscia 281
storia e architettura 112
storia e caratteristica 105
storia e cultura 85
storia e descrizione 138
storia e filosofia 80
storia e leggenda 138
storia e profilo 128
strada e autostrada 66
strada e piazza 63
strada ed autostrada 76
struttura e funzione 56
studente e docente 55
studio e ricerca 298
tema e documento 126
temperatura e pressione 217
tempo e denaro 101
tempo e luogo 82
tempo e modo 86
tempo e spazio 71
teologia e filosofia 65
teoria e pratica 76
testo e disegno 78
testo e musica 158
tetta e culo 59
traccia e formato 61
trasporto e via 72
trasporto e viabilità 84
trasporto ed infrastruttura 62
uomo e animale 89
uomo e donna 2243
uomo e mezzo 285
uomo o donna 114
uso e consumo 216
uso e costume 348
vantaggio e svantaggio 174
varietà e sottospecie 100
venerdì e sabato 99
vescovo e arcivescovo 403
via e piazza 76
viaggio e reportage 120
video e audio 175

villa e palazzo 57	vitto e alloggio 99
vincitore e candidato 213	vittoria e sconfitta 58
violino e orchestra 66	voce e chitarra 205
violino e pianoforte 89	voce e pianoforte 56
vita e carriera 101	volta e mezzo 62
vita e morte 125	volume e mezzo 58
vita e opera 220	
vitamina e minerale 85	

Di seguito viene inoltre riportato l'insieme dei dati prodotti dallo strumento computazionale per tutte le espressioni considerate. I nomi delle colonne sono precisati dalla seguente legenda:

ESPRESSIONE	Espressione considerata
FREQ	Frequenza di occorrenza (componenti lemmatizzati)
FREQ_INT1	Numero di occorrenze dell'espressione interrotta tra primo e secondo componente
FREQ_INT2	Numero di occorrenze dell'espressione interrotta tra secondo e terzo componente
INTERROMP	Valore di I_{syn}^{int}
FREQ_REV	Numero di occorrenze dell'espressione con componenti pieni invertiti
ORDINE INVERSO	Valore di I_{syn}^{ord}
MOD_SINTAG	Valore di I_{syn}
FREQ_SIN1	Numero di occorrenze delle espressioni con primo componente sostituito (valore non intero poiché normalizzato in base alla similarità semantica)
DISP1	Dispersione della sostituzione del primo componente
FREQ_SIN2	Numero di occorrenze delle espressioni con secondo componente pieno sostituito (valore non intero poiché normalizzato in base alla similarità semantica)
DISP2	Dispersione della sostituzione del secondo componente
MOD_PARAD	Valore di I_{sub}
FORMA_PREVAL	Forma con il più alto numero di occorrenze
FREQ_P	Numero di occorrenze della forma prevalente
FORMA_PREVAL_2	Forma con il secondo numero di più alto di occorrenze
FREQ_P_2	Numero di occorrenze della seconda forma prevalente
MOD_FLESS	Valore di I_{infl}
ERR	Se contrassegnato da X, espressione scartata per errori dovuti al corpus

I

Dati sul pattern VCV

Si riporta di seguito la lista, in ordine alfabetico, delle 500 espressioni relative al pattern VCV estratte in PAISÀ e i dati relativi al comportamento variazionale di ognuna delle espressioni.

abbandonare e sostituire 26	ampliare e trasformare 20
abbattere e ricostruire 35	ampliare ed abbellire 19
abbattere e sostituire 19	analizzare e studiare 24
abolire e sostituire 27	andare e tornare 40
accendere e spegnere 29	andare e venire 144
accettare o rifiutare 37	apparire e scomparire 57
addestrare ed equipaggiare 32	aprire e chiudere 186
affrontare e battere 26	aprire o chiudere 58
affrontare e risolvere 82	arrestare e accusare 27
affrontare e sconfiggere 117	arrestare e condannare 247
affrontare e superare 33	arrestare e condurre 49
affrontare e uccidere 19	arrestare e deportare 50
affrontare e vincere 24	arrestare e giustiziare 35
aggiungere o rimuovere 20	arrestare e imprigionare 58
aggiungere o togliere 24	arrestare e incarcerare 55
aggredire e uccidere 21	arrestare e mandare 21
aiutare e sostenere 22	arrestare e mettere 22
allevare e riprodurre 63	arrestare e portare 68
alzare e andare 29	arrestare e processare 79
alzare o abbassare 28	arrestare e rilasciare 19
amare e odiare 33	arrestare e rinchiudere 51
amare e rispettare 47	arrestare e torturare 38
amare o odiare 21	arrestare e tradurre 19
ammalare e morire 84	arrestare ed imprigionare 32
ampliare e migliorare 22	arrestare ed incarcerare 20
ampliare e modificare 29	arsenati e vanadati 30
ampliare e restaurare 22	aspettare e vedere 32
ampliare e ristrutturare 29	assediare e conquistare 52

assediare e distruggere 31
 assediare e prendere 19
 attaccare e conquistare 52
 attaccare e distruggere 97
 attaccare e saccheggiare 35
 attaccare e sconfiggere 41
 attaccare e uccidere 43
 attaccare ed uccidere 24
 attivare o disattivare 31
 aumentare e diminuire 21
 aumentare o diminuire 123
 avere e avere 28
 avere ed avere 30
 ballare e cantare 68
 bastare e avanzare 86
 bastare ed avanzare 20
 battere e ribattere 22
 bisognare però ricordare 21
 cancellare e sostituire 20
 cantare e ballare 110
 cantare e danzare 23
 cantare e recitare 32
 cantare e suonare 124
 caricare e scaricare 28
 catturare e condannare 43
 catturare e condurre 25
 catturare e decapitare 19
 catturare e fare 34
 catturare e fucilare 26
 catturare e giustiziare 67
 catturare e imprigionare 57
 catturare e mettere 26
 catturare e portare 85
 catturare e rinchiudere 35
 catturare e tenere 26
 catturare e torturare 49
 catturare e uccidere 77
 catturare ed imprigionare 43
 catturare ed uccidere 39
 catturare o uccidere 31
 censire e documentare 38
 cercare e trovare 75
 chiedere e chiedere 19

chiedere e ottenere 263
 chiedere ed ottenere 364
 citare o riprodurre 25
 collaborare e collaborare 40
 colpire e distruggere 26
 colpire e uccidere 32
 combattere e morire 19
 combattere e sconfiggere 35
 combattere e vincere 54
 comporre e cantare 47
 comporre e interpretare 34
 comporre e produrre 21
 comporre e pubblicare 20
 comporre e registrare 29
 comporre e suonare 28
 comporre ed eseguire 38
 comporre ed interpretare 54
 comprare e vendere 53
 comprare o vendere 19
 comprendere ed individuare 48
 concepire e partorire 23
 concepire e realizzare 43
 condannare e giustiziare 23
 condannare o uccidere 30
 confermare o smentire 36
 conoscere e amare 34
 conoscere e apprezzare 97
 conoscere e diffondere 19
 conoscere e divenire 28
 conoscere e fare 26
 conoscere e frequentare 86
 conoscere e rispettare 23
 conoscere e sposare 83
 conoscere e stimare 42
 conoscere e studiare 23
 conoscere e utilizzare 32
 conoscere ed apprezzare 123
 conquistare e distruggere 36
 conquistare e saccheggiare 34
 conservare ed esporre 49
 continuare e continuare 21
 controllare e garantire 110
 controllare e gestire 35

-
- convocare e presiedere 26
copiare e incollare 36
copiare ed incollare 23
correre e vincere 23
costruire e gestire 37
costruire e mantenere 23
creare e dirigere 48
creare e gestire 45
creare e mantenere 24
creare e produrre 36
creare e pubblicare 23
creare e sviluppare 29
credere e sperare 25
crescere e bastare 255
crescere e divenire 23
crescere e diventare 49
crescere e maturare 28
crescere e morire 20
crescere ed educare 21
custodire ed esporre 19
danneggiare o distruggere 46
dare e ricevere 33
datare e firmare 25
decollare e atterrare 27
decollare ed atterrare 19
demolire e ricostruire 78
demolire e sostituire 31
deporre ed esiliare 20
devastare e saccheggiare 21
dimettere e subentrare 25
dire e fare 98
dire e ripetere 28
dire e scrivere 59
dire o fare 56
dire o scrivere 22
dirigere e interpretare 202
dirigere e produrre 101
dirigere e sceneggiare 47
dirigere ed interpretare 99
disegnare e costruire 33
disegnare e dipingere 19
disegnare e realizzare 35
disputare e vincere 26
distruggere e ricostruire 73
distruggere e saccheggiare 27
distruggere o danneggiare 42
eliminare e sostituire 23
entrare e uscire 153
entrare ed uscire 209
entrare o uscire 49
esistere e essere 34
esistere ed esistere 20
esistere ed essere 63
esistere o essere 19
esonerare e sostituire 46
essere e continuare 56
essere e essere 280
essere e fare 21
essere e restare 173
essere e rimanere 167
essere ed essere 471
essere o essere 53
estrarre e lavorare 19
fare e bastare 20
fare e continuare 28
fare e dire 48
fare e disfare 39
fare e essere 25
fare e fare 75
fare ed essere 22
fare o dire 41
ferire e fare 28
ferire o uccidere 25
fermare e arrestare 19
firmare e datare 168
fondare e dirigere 411
fondare e gestire 26
fondare e guidare 40
giocare e vincere 44
girare e rigirare 33
giudicare e condannare 45
grattare e vincere 24
guardare e passare 21
ideare e condurre 130
ideare e costruire 33
ideare e curare 27

- ideare e dirigere 70
 ideare e disegnare 28
 ideare e produrre 61
 ideare e progettare 22
 ideare e realizzare 118
 ideare e scrivere 32
 ideare e sviluppare 23
 ideare ed organizzare 21
 importare ed esportare 19
 imprigionare e torturare 29
 inalare o ingerire 21
 incendiare e distruggere 26
 inciampare e cadere 21
 incidere e pubblicare 21
 incontrare e sposare 33
 indagare e scoprire 21
 iniziare e finire 59
 iniziare e terminare 48
 innamorare e sposare 27
 inseguire e uccidere 19
 intagliare e dipingere 19
 intagliare e dorare 25
 intendere e volere 50
 interpretare e dirigere 54
 interrogare e torturare 22
 invadere e conquistare 33
 invadere e distruggere 21
 invadere e occupare 28
 invadere e saccheggiare 34
 inventare e brevettare 23
 investire e uccidere 21
 inviare e ricevere 48
 kmò ed avere 0
 kmò ed essere 0
 lavorare e produrre 27
 lavorare e vivere 30
 lavorare o lavorare 53
 legare e imbavagliare 28
 leggere e capire 20
 leggere e commentare 28
 leggere e rileggere 43
 leggere e scrivere 377
 leggere e studiare 21
 leggere né scrivere 21
 leggere o scrivere 43
 mangiare e bere 139
 mangiare e dormire 29
 minare e fare 19
 modificare e ampliare 19
 modificare e rendere 19
 molibdati e tungstati 34
 mordere e fuggire 73
 morire e essere 85
 morire e ferire 23
 morire e risorgere 30
 morire e sepolto 29
 nascere e crescere 549
 nascere e morire 83
 nascere e sviluppare 43
 nascere e vivere 146
 nominare e revocare 19
 occupare e saccheggiare 19
 organizzare e condurre 20
 organizzare e dirigere 58
 organizzare e gestire 64
 organizzare e partecipare 26
 organizzare e promuovere 28
 osservare e descrivere 22
 parere e piacere 23
 parlare e scrivere 65
 parlare o scrivere 27
 partecipare e vincere 96
 partire e arrivare 26
 pensare e agire 22
 pensare e realizzare 47
 piangere e litigare 214
 possedere e gestire 39
 precedere e seguire 75
 precedere o seguire 28
 prendere e distruggere 19
 prendere e saccheggiare 31
 prendere o lasciare 109
 processare e condannare 190
 processare e giustiziare 25
 produrre e arrangiare 41
 produrre e commercializzare 87

-
- produrre e consumare 56
produrre e dirigere 111
produrre e distribuire 175
produrre e interpretare 21
produrre e mettere 20
produrre e pubblicare 21
produrre e realizzare 38
produrre e scrivere 30
produrre e trasmettere 69
produrre e vendere 125
produrre ed esportare 23
produrre ed interpretare 25
progettare e costruire 467
progettare e dirigere 24
progettare e produrre 87
progettare e realizzare 380
progettare e sviluppare 39
promuovere e coordinare 24
promuovere e diffondere 21
promuovere e organizzare 27
promuovere e sostenere 50
promuovere e sviluppare 22
promuovere ed organizzare 22
provare e riprovare 35
pubblicare e distribuire 36
puntare e clicca 19
raccolgere e catalogare 23
raccolgere e conservare 41
raccolgere e pubblicare 85
raccolgere ed elaborare 20
raggiungere e superare 196
raggiungere o superare 50
rapire e portare 53
rapire e torturare 26
rapire e uccidere 58
realizzare e gestire 26
realizzare e produrre 20
realizzare e pubblicare 19
recitare e cantare 23
recitare o cantare 25
registrare e mixato 43
registrare e pubblicare 85
registrare e trasmettere 20
restaurare e ampliare 20
restaurare e riaprire 35
restaurare e riportare 27
restaurare ed ampliare 42
richiedere e ottenere 21
richiedere ed ottenere 35
riconoscere e accettare 19
riconoscere e garantire 39
riconoscere e rispettare 24
riconoscere e tutelare 23
ricostruire e ampliare 25
ridere e fare 20
ridere e scherzare 48
ridere o piangere 25
ridurre o eliminare 29
rimuovere e sostituire 45
rinsaldare e unificare 19
riprendere e ampliare 20
riprendere e sviluppare 70
ristrutturare e ampliare 28
ristrutturare e trasformare 19
ristrutturare ed ampliare 39
rivedere e correggere 19
riveduto e correggere 28
saccheggiare e bruciare 27
saccheggiare e dare 20
saccheggiare e devastare 23
saccheggiare e distruggere 65
saccheggiare e incendiare 36
saccheggiare ed incendiare 26
salare e pepare 28
salire e scendere 95
salire o scendere 39
sapere né leggere 26
scaricare e installare 20
scaricare ed installare 22
sceneggiare e dirigere 20
sconfiggere e catturare 51
sconfiggere e costringere 23
sconfiggere e fare 30
sconfiggere e uccidere 196
sconfiggere ed uccidere 45
scoprire e denunciare 20

- scoprire e descrivere 22
scoprire e lanciare 54
scoprire e mettere 19
scoprire e uccidere 21
scrivere e arrangiare 21
scrivere e cantare 149
scrivere e comporre 77
scrivere e condurre 55
scrivere e dire 19
scrivere e dirigere 851
scrivere e disegnare 398
scrivere e fare 38
scrivere e firmare 21
scrivere e ideare 24
scrivere e illustrare 39
scrivere e incidere 21
scrivere e interpretare 128
scrivere e leggere 38
scrivere e mettere 22
scrivere e musicare 33
scrivere e parlare 61
scrivere e produrre 262
scrivere e pubblicare 212
scrivere e realizzare 25
scrivere e recitare 22
scrivere e registrare 120
scrivere e suonare 40
scrivere ed eseguire 31
scrivere ed ideare 35
scrivere ed interpretare 136
sedurre e abbandonare 24
segnalare o recensire 24
sentire e vedere 26
smontare e rimontare 33
sopprimere e aggregare 27
sopprimere e sostituire 39
sorgere e tramontare 28
sostenere e finanziare 20
sostenere e promuovere 30
sparare e uccidere 48
sparare ed uccidere 21
sposare e avere 123
sposare e divorziare 24
sposare e vivere 21
sposare ed avere 150
staccare e conservare 24
stringere e allungare 26
studiare e fare 22
studiare e lavorare 27
studiare e realizzare 28
suonare e cantare 111
suonare e registrare 22
sviluppare e costruire 34
sviluppare e diffondere 21
sviluppare e distribuire 29
sviluppare e mantenere 21
sviluppare e produrre 67
sviluppare e pubblicare 52
sviluppare ed appuntire 31
svolgere e svolgere 23
tagliare e cucire 27
temere e rispettare 35
tirare e mollare 94
toccare e superare 19
torturare e uccidere 70
torturare ed uccidere 37
tradurre e adattare 31
tradurre e commentare 22
tradurre e pubblicare 133
trasmettere e ricevere 25
travolgere e uccidere 19
trovare e distruggere 33
trovare e uccidere 20
tutelare e valorizzare 21
uccidere e fare 19
uccidere o catturare 33
uccidere o fare 21
uccidere o ferire 38
usare e gettare 149
uscire e andare 27
validare e standardizzare 19
vedere e ascoltare 23
vedere e commentare 92
vedere e conoscere 19
vedere e considerare 99
vedere e rivedere 75

vedere e sentire 147
vedere e toccare 23
vedere e vivere 29
vedere o sentire 33
vendere e comprare 19
vincere e convincere 20
vincere o perdere 42
violentare e uccidere 32

vivere e lavorare 690
vivere e morire 97
vivere e operare 44
vivere e studiare 27
vivere ed operare 38
vivere o morire 39

Di seguito viene inoltre riportato l'insieme dei dati prodotti dallo strumento computazionale per tutte le espressioni considerate. I nomi delle colonne sono precisati dalla seguente legenda:

ESPRESSIONE	Espressione considerata
FREQ	Frequenza di occorrenza (componenti lemmatizzati)
FREQ_INT1	Numero di occorrenze dell'espressione interrotta tra primo e secondo componente
FREQ_INT2	Numero di occorrenze dell'espressione interrotta tra secondo e terzo componente
INTERROMP	Valore di I_{syn}^{int}
FREQ_REV	Numero di occorrenze dell'espressione con componenti pieni invertiti
ORDINE INVERSO	Valore di I_{syn}^{ord}
MOD_SINTAG	Valore di I_{syn}
FREQ_SIN1	Numero di occorrenze delle espressioni con primo componente sostituito (valore non intero poiché normalizzato in base alla similarità semantica)
DISP1	Dispersione della sostituzione del primo componente
FREQ_SIN2	Numero di occorrenze delle espressioni con secondo componente pieno sostituito (valore non intero poiché normalizzato in base alla similarità semantica)
DISP2	Dispersione della sostituzione del secondo componente
MOD_PARAD	Valore di I_{sub}
FORMA_PREVAL	Forma con il più alto numero di occorrenze
FREQ_P	Numero di occorrenze della forma prevalente
FORMA_PREVAL_2	Forma con il secondo numero di più alto di occorrenze
FREQ_P_2	Numero di occorrenze della seconda forma prevalente
MOD_FLESS	Valore di I_{infl}
ERR	Se contrassegnato da X, espressione scartata per errori dovuti al corpus

[illegible]

[illegible]

J

Dati sul pattern VDN

Si riporta di seguito la lista, in ordine alfabetico, delle 500 espressioni relative al pattern VDN estratte in PAISÀ e i dati relativi al comportamento variazionale di ognuna delle espressioni.

abbandonare il città 357
abbandonare il gruppo 509
abbandonare il progetto 373
abbandonare il studio 427
accettare il invito 343
accettare il proposta 608
acquisire il diritto 359
acquistare il diritto 440
affidare il compito 500
affidare il incarico 365
affrontare il problema 713
affrontare il tema 548
andare il cosa 368
aprire il fuoco 709
aprire il occhio 539
aprire il porta 1583
aprire il strada 1045
arricchire il lista 426
arrivare il momento 750
assumere il aspetto 328
assumere il carica 456
assumere il comando 695
assumere il controllo 470
assumere il denominazione 739
assumere il direzione 469
assumere il forma 627
assumere il incarico 505
assumere il nome 2165

assumere il ruolo 741
assumere il titolo 525
attirare il attenzione 1497
attraversare il città 353
attraversare il fiume 528
attraversare il territorio 344
aumentare il numero 581
avanzare il ipotesi 377
avere il aspetto 477
avere il capacità 1026
avere il capello 417
avere il certezza 326
avere il compito 2452
avere il controllo 478
avere il coraggio 1530
avere il diritto 1693
avere il dovere 436
avere il effetto 451
avere il forma 885
avere il fortuna 706
avere il forza 492
avere il funzione 920
avere il idea 693
avere il impressione 833
avere il incarico 470
avere il meglio 1925
avere il merito 512
avere il nome 508

- avere il obbligo 395
avere il obiettivo 381
avere il occasione 647
avere il occhio 435
avere il onore 531
avere il opportunità 673
avere il possibilità 3334
avere il potere 1135
avere il ruolo 611
avere il scopo 1972
avere il sensazione 378
avere il tempo 754
avere il titolo 495
avere il vantaggio 476
cambiare il mondo 566
cambiare il nome 1037
catturare il attenzione 335
causare il morte 839
cedere il passo 338
cessare il fuoco 526
chiedere il aiuto 369
chiudere il battente 353
chiudere il carriera 475
chiudere il occhio 398
chiudere il porta 358
cogliere il occasione 1120
collegare il città 448
compiere il studio 484
completare il studio 495
comprendere il città 417
comprendere il parte 328
coniare il termine 360
conoscere il nome 364
conquistare il città 435
conquistare il promozione 422
conquistare il titolo 688
conseguire il diploma 497
conseguire il dottorato 331
conseguire il laurea 743
consentire il accesso 337
considerare il padre 325
contattare il concessionario 1216
continuare il attività 336
continuare il studio 430
correre il rischio 612
costare il vita 439
costituire il base 482
creare il condizione 354
dare il caccia 707
dare il colpa 396
dare il dimissione 412
dare il idea 459
dare il impressione 619
dare il meglio 372
dare il natale 691
dare il nome 2594
dare il ordine 368
dare il possibilità 1473
dare il titolo 618
dare il via 2108
deporre il uovo 390
derivare il nome 702
detenere il record 362
dire il nome 344
dire il verità 1098
disputare il campionato 411
durare il giornale 344
effettuare il login 5518
esercitare il professione 492
essere il acronimo 396
essere il album 465
essere il amore 329
essere il anno 1865
essere il area 325
essere il attività 389
essere il attore 402
essere il autore 1013
essere il base 996
essere il canzone 373
essere il capacità 478
essere il capitale 1012
essere il capitano 395
essere il capo 957
essere il capoluogo 457
essere il casa 562
essere il caso 3394

-
- essere il causa 1266
essere il centro 1294
essere il chiave 392
essere il chiesa 1271
essere il città 2554
essere il colore 359
essere il comandante 322
essere il condizione 462
essere il conseguenza 400
essere il cosa 675
essere il costruzione 358
essere il cuore 347
essere il differenza 356
essere il dimostrazione 333
essere il direttore 593
essere il donna 477
essere il effetto 347
essere il elemento 885
essere il erede 363
essere il esempio 381
essere il espressione 381
essere il famiglia 403
essere il fatto 1557
essere il figlio 2760
essere il figura 406
essere il film 514
essere il fine 728
essere il fondatore 634
essere il fonte 427
essere il forma 794
essere il forza 468
essere il fratello 897
essere il frutto 970
essere il funzione 433
essere il giorno 545
essere il gruppo 614
essere il guerra 329
essere il idea 370
essere il inizio 1048
essere il insieme 731
essere il leader 464
essere il legge 345
essere il linea 337
essere il lingua 640
essere il luogo 1347
essere il madre 578
essere il massimo 350
essere il metodo 384
essere il mezzo 364
essere il modello 585
essere il modo 764
essere il moglie 439
essere il momento 1610
essere il mondo 327
essere il motivo 732
essere il motore 372
essere il necessità 363
essere il nome 5784
essere il numero 1308
essere il occasione 593
essere il oggetto 364
essere il opera 622
essere il organo 430
essere il padre 1149
essere il paese 505
essere il parola 553
essere il parte 1214
essere il periodo 584
essere il persona 672
essere il personaggio 641
essere il posizione 326
essere il possibilità 1246
essere il presenza 1377
essere il presidente 574
essere il prima 654
essere il principe 344
essere il problema 658
essere il processo 354
essere il prodotto 663
essere il produzione 425
essere il progetto 351
essere il programma 350
essere il proprietario 344
essere il protagonista 1311
essere il prova 586
essere il punto 1502

essere il ragazzo 394
essere il ragione 427
essere il rapporto 507
essere il rappresentazione 337
essere il re 385
essere il responsabile 603
essere il resto 336
essere il rischio 469
essere il risposta 501
essere il risultato 2369
essere il ruolo 663
essere il scelta 439
essere il secondo 642
essere il sede 1249
essere il segno 338
essere il serie 345
essere il simbolo 775
essere il sistema 709
essere il situazione 337
essere il società 423
essere il soluzione 464
essere il somma 321
essere il sorella 469
essere il squadra 626
essere il stato 487
essere il stazione 468
essere il storia 1148
essere il strada 449
essere il strumento 460
essere il studio 427
essere il tema 405
essere il tempo 696
essere il tentativo 339
essere il termine 713
essere il tipo 505
essere il titolo 2096
essere il turno 477
essere il ultimo 451
essere il unico 1298
essere il uomo 962
essere il uso 595
essere il valore 454
essere il velocità 403

essere il verità 352
essere il versione 1237
essere il via 376
essere il vincitore 361
essere il vita 620
essere il voce 524
essere il volontà 324
essere il volta 1915
essere il zona 427
evolvere il pubblicità 344
fare il amore 836
fare il bagno 347
fare il conoscenza 628
fare il conte 1943
fare il cosa 520
fare il differenza 604
fare il giro 623
fare il guerra 323
fare il nome 490
fare il punto 380
fare il spesa 486
fare il storia 519
favorire il sviluppo 444
fondare il città 404
frequentare il corso 386
frequentare il scuola 1123
garantire il sicurezza 347
gettare il base 612
giocare il partita 418
girare il film 393
girare il mondo 424
giungere il momento 583
indicare il numero 329
indossare il maglia 464
informare il giovane 348
iniziare il attività 602
iniziare il carriera 1393
iniziare il costruzione 769
iniziare il lavoro 830
iniziare il produzione 444
iniziare il studio 574
interpretare il parte 539
interpretare il personaggio 393

-
- interpretare il ruolo 2320
intraprendere il carriera 1361
introdurre il concetto 368
lasciare il band 998
lasciare il casa 343
lasciare il città 657
lasciare il gruppo 1236
lasciare il paese 362
lasciare il posto 1121
lasciare il scuola 373
lasciare il segno 430
leggere il libro 371
mantenere il controllo 469
mettere il mano 850
migliorare il condizione 336
migliorare il prestazione 354
migliorare il qualità 387
narrare il storia 741
narrare il vicenda 544
occupare il città 334
offrire il possibilità 703
ottenere il cattedra 329
ottenere il indipendenza 364
ottenere il permesso 446
ottenere il promozione 476
ottenere il riconoscimento 373
ottenere il risultato 351
ottenere il titolo 739
pagare il cauzione 593
pagare il giornale 347
pagare il tassa 529
passare il notte 582
passare il tempo 537
percorrere il strada 378
perdere il controllo 846
perdere il filo 458
perdere il lavoro 342
perdere il senso 395
perdere il testa 354
perdere il traccia 448
perdere il vita 2034
permettere il accesso 339
permettere il passaggio 321
porre il accento 485
porre il base 743
porre il problema 658
portare il nome 1089
portare il squadra 349
prendere il comando 679
prendere il controllo 1337
prendere il decisione 847
prendere il distanza 702
prendere il forma 337
prendere il nome 10445
prendere il parola 359
prendere il posto 2251
prendere il potere 911
prendere il sopravvento 640
prendere il via 1337
prendere il voto 357
prestare il voce 380
prevedere il costruzione 354
prevedere il possibilità 392
prevedere il uso 363
prevedere il utilizzo 407
proseguire il studio 517
provocare il morte 590
pubblicare il album 885
pubblicare il libro 432
puntare il dito 427
raccontare il storia 1632
raccontare il vicenda 347
raccontare il vita 356
raggiungere il apice 561
raggiungere il età 393
raggiungere il finale 625
raggiungere il grado 343
raggiungere il livello 355
raggiungere il maturità 367
raggiungere il numero 413
raggiungere il obiettivo 463
raggiungere il posizione 1130
raggiungere il semifinale 361
raggiungere il successo 397
raggiungere il velocità 393
raggiungere il vetta 760

- rassegnare il dimissione 614
recitare il parte 357
restare il fatto 390
ricevere il incarico 413
ricevere il ordine 564
ricevere il premio 781
ricevere il titolo 540
ricevere il visita 634
richiamare il attenzione 448
richiedere il rimozione 2761
riconoscere il diritto 357
ricoprire il carica 1155
ricoprire il incarico 645
ricoprire il ruolo 1424
ricordare il esperienza 347
ridurre il costo 446
ridurre il numero 436
riguardare il composizione 558
riguardare il suddivisione 525
rinnovare il contratto 365
riportare il nome 364
riprendere il attività 347
riprendere il controllo 403
rischiare il vita 405
riservare il diritto 577
risolvere il problema 2480
risolvere il questione 417
salvare il informazione 6650
salvare il mondo 354
salvare il vita 1570
scoppiare il guerra 390
scoprire il verità 544
scrivere il sceneggiatura 351
scrivere il testo 387
segnalare il presenza 378
segnare il confine 643
segnare il fine 773
segnare il gol 414
segnare il inizio 953
segnare il ritorno 330
seguire il consiglio 323
seguire il corso 551
seguire il esempio 448
seguire il indicazione 345
seguire il orma 655
seguire il sorta 410
seguire il vicenda 416
sentire il bisogno 508
sentire il necessità 342
soddisfare il esigenza 335
sorgere il chiesa 325
sostenere il necessità 327
sottolineare il importanza 324
spianare il strada 326
sposare il figlio 520
stabilire il record 329
stare il cosa 399
successe il figlio 400
suonare il batteria 343
suonare il chitarra 973
suonare il pianoforte 347
superare il limite 474
svolgere il attività 415
svolgere il funzione 690
svolgere il ruolo 590
tentare il carrambata 452
tentare il suicidio 389
terminare il studio 649
togliere il vita 390
trovare il cattedrale 1225
trovare il chiesa 334
trovare il modo 711
trovare il morte 704
usare il nome 362
usare il parola 369
usare il termine 808
uscire il album 416
utilizzare il termine 407
valere il pena 1765
vedere il film 492
vedere il luce 1645
vedere il ora 926
vedere il partecipazione 1074
vedere il presenza 450
vedere il successo 373
vedere il vittoria 658

vestire il maglia 931	vincere il premio 1765
vestire il panno 823	vincere il scudetto 392
vincere il battaglia 326	vincere il serata 743
vincere il campionato 2369	vincere il titolo 1498
vincere il concorso 622	vincere il torneo 507
vincere il elezione 761	visualizzare il commento 481
vincere il gara 521	
vincere il medaglia 986	

Di seguito viene inoltre riportato l'insieme dei dati prodotti dallo strumento computazionale per tutte le espressioni considerate. I nomi delle colonne sono precisati dalla seguente legenda:

ESPRESSIONE	Espressione considerata
FREQ	Frequenza di occorrenza (componenti lemmatizzati)
F_INT1	Numero di occorrenze dell'espressione interrotta tra primo e secondo componente
F_INT2	Numero di occorrenze dell'espressione interrotta tra secondo e terzo componente
INTERROMP	Valore di I_{syn}^{int}
F_T1	Numero di occorrenze dell'espressione topicalizzata (verbo semplice)
F_T2	Numero di occorrenze dell'espressione topicalizzata (verbo con ausiliare)
TOPICALIZZ	Valore dell'indice di topicalizzazione
F_A1	Numero di occorrenze con ripresa anaforica (verbo semplice)
F_A2	Numero di occorrenze con ripresa anaforica (verbo con ausiliare)
ANAFORA	Valore dell'indice di anaforizzazione
F_P1	Numero di occorrenze con trasformazione passiva (un verbo ausiliare)
F_A2	Numero di occorrenze con trasformazione passiva (due verbi ausiliari)
PASSIVO	Valore dell'indice di passivizzazione
F_R1	Numero di occorrenze con trasformazione relativa (verbo semplice)
F_R2	Numero di occorrenze con trasformazione relativa (un verbo ausiliare)
F_R3	Numero di occorrenze con trasformazione relativa (due verbi ausiliari)
RELATIVO	Valore dell'indice di relativizzazione
MOD_SINTAG	Valore di I_{syn}
FREQ_SIN1	Numero di occorrenze delle espressioni con primo componente sostituito (valore non intero poiché normalizzato in base alla similarità semantica)
DISP1	Dispersione della sostituzione del primo componente
FREQ_SIN2	Numero di occorrenze delle espressioni con secondo componente pieno sostituito (valore non intero poiché normalizzato in base alla similarità semantica)
DISP2	Dispersione della sostituzione del secondo componente
MOD_PARAD	Valore di I_{sub}
FORMA_PREVAL	Forma con il più alto numero di occorrenze
FREQ_P	Numero di occorrenze della forma prevalente
FORMA_PREVAL_2	Forma con il secondo numero di più alto di occorrenze
FREQ_P_2	Numero di occorrenze della seconda forma prevalente
MOD_FLESS	Valore di I_{infl}
ERR	Se contrassegnato da X, espressione scartata per errori dovuti al corpus

[illegible]

[illegible]

K

Part of speech Tagset - Fisica

Tagset dell'italiano in uso nel *parameter file* di TreeTagger (Schmid, 1994) -
Copyright Prof. Achim Stein, University of Stuttgart.

ABR	abbreviation
ADJ	adjective
ADV	adverb
CON	conjunction
DET:def	definite article
DET:indef	indefinite article
FW	foreign word
INT	interjection
LS	list symbol
NOM	noun
NPR	name
NUM	numeral
PON	punctuation
PRE	preposition
PRE:det	preposition+article
PRO	pronoun
PRO:demo	demonstrative pronoun
PRO:indef	indefinite pronoun
PRO:inter	interrogative pronoun
PRO:pers	personal pronoun
PRO:poss	possessive pronoun
PRO:refl	reflexive pronoun
PRO:rela	relative pronoun
SENT	sentence marker
SYM	symbol
VER:cimp	verb conjunctive imperfect
VER:cond	verb conditional
VER:cpre	verb conjunctive present
VER:futu	verb future tense
VER:geru	verb gerund
VER:impe	verb imperative
VER:impf	verb imperfect
VER:infi	verb infinitive
VER:pper	verb participle perfect
VER:ppre	verb participle present
VER:pres	verb present
VER:refl:infi	verb reflexive infinitive
VER:remo	verb simple past

Bibliografia

- ALISOVA T. A. 1967. Studi di sintassi italiana. *Studi di Filologia Italiana*, **XXV**, 223–313.
- ATTARDI G., DELL’ORLETTA F., SIMI M., & TURIAN J. 2009. Accurate Dependency Parsing with a Stacked Multilayer Perceptron. *In: Proceedings of Evalita ‘09, Evaluation of NLP and Speech Tools for Italian*.
- BALDWIN T., BANNARD C., TANAKA T., & WIDDOWS D. 2003. An empirical model of multiword expression decomposability. *Pages 89–96 of: Proceedings of ACL-SIGLEX Workshop on Multiword Expressions*.
- BALLY C. 1951. *Traité de stylistique française*. Paris, Klincksieck. I edn. 1909.
- BANNARD C. 2007. A measure of syntactic flexibility for automatically identifying multiword expressions in corpora. *In: Proceedings of ACL 2007 (Workshop)*.
- BANNARD C., BALDWIN T., & LASCARIDES A. 2003. A statistical approach to the semantics of verb-particles. *Pages 65–72 of: Proceedings of ACL-SIGLEX Workshop on Multiword Expressions*.
- BARONI M., GUEVARA E., PIRRELLI V., & ZANCHETTA E. 2006. Corpus evidence and compound structure: the case of Italian NN compounds. *In: Proceedings of Quantitative Investigations in Theoretical Linguistics 2 (QITL-2)*.
- BARRACHINA S., BENDER O., CASACUBERTA F., CIVERA J., CUBEL E., KHADIVI S., LAGARDA A., NEY H., TOMÁS J., VIDAL E., & VILAR J.-M. 2009. Statistical Approaches to Computer-Assisted Translation. *Computational Linguistics*, **35**(1), 3–28.
- BARREIRO A. 2008. *Make it simple with paraphrases. Automated paraphrasing for authoring aids and machine translation*. Ph.D. thesis, Faculdade de Letras da Universidade do Porto, Oporto, Portugal.
- BENVENISTE E. 1966. Différentes formes de la composition nominale en français. *Bulletin de la Société de Linguistique de Paris*, **61**(1), 82–95.
- BLOOMFIELD L. 1967. *Language*. London, George Allen & Unwin LTD. I ed. 1933.
- BNC. 2001. *The British National Corpus*. Distributed by Oxford University Computing Services on behalf of the BNC Consortium. <http://www.natcorp.ox.ac.uk/>.

- BOLASCO S. 1999. *Analisi multidimensionale dei dati. Metodi, strategie e criteri d'interpretazione*. Roma: Carocci.
- BOLASCO S. 2013. *L'analisi automatica dei testi. Fare ricerca con il text mining*. Roma: Carocci.
- BOSQUE I. 2004a. Combinatoria y significación. Algunas reflexiones. In: *REDES, Diccionario Combinatorio del Español Contemporáneo*. Hoepli.
- BOSQUE I. 2004b. *REDES, Diccionario Combinatorio del Español Contemporáneo*. Hoepli.
- BOWKER L. 2002. *Computer-Aided Translation Technology*. University of Ottawa Press.
- BRÉAL M. 1904. *Essai de sémantique*. III edn. Paris, Hachette. I edizione 1897.
- BRAME M. 1984a. Ungrammatical notes 5: asymmetry and the creation. *Linguistic Analysis*, 51–54.
- BRAME M. 1984b. Universal word induction vs move a. *Linguistic Analysis*, 313–352.
- BREIMAN L. 2001. Statistical Modeling: The Two Cultures. *Statistical Science*, 16(3), 199–231.
- CALZOLARI N., FILLMORE C., GRISHAM R., IDE N., LENCI A., MCLEOD C., & ZAMPOLLI A. 2002. Towards best practise for multiword expressions in computational lexicons. In: *Proceedings of LREC 2002: 3rd International Conference on Language Resources and Evaluation*.
- CAP F., WELLER M., & HEID U. 2013. Using a Rich Feature Set for the Identification of German MWEs. In: *Proceedings of Machine Translation Summit XIV, Nice, France*.
- CASADEI F. 1994. Il lessico nelle strategie di presentazione dell'informazione scientifica: il caso della fisica. Pages 47–69 of: DE MAURO T. (ed), *Studi sul trattamento linguistico dell'informazione scientifica*. Roma: Bulzoni.
- CASADEI F. 1996. *Metafore ed espressioni idiomatiche. Uno studio semantico sull'italiano*. Bulzoni Editore.
- CHAFE W. L. 1968. Idiomaticity as an anomaly in the Chomskyan paradigm. *Foundations of Language*, 109–127.
- CHIARI I. 2007. *Introduzione alla linguistica computazionale*. Laterza.
- CHOMSKY N. 1957. *Syntactic Structures*. Mouton & co.

- CHOMSKY N. 1980. *Rules and representation*. New York, Columbia University Press.
- CHOUKA Y. 1988. Looking for Needles in a Haystack. *Pages 609–623 of: Proceedings of RIAO '88*.
- CHURCH K. W., & HANKS P. 1990. Word association norms, mutual information, and lexicography. *Computational Linguistics*, 1(16), 22–29.
- CHURCH K. W., GALE W. A., HANKS P., & HINDLE D. 1991. Using statistics in lexical analysis. In: ZERNIK U. (ed), *Lexical Acquisition: Exploiting On-line Resources to Build a Lexicon*. Hillsdale, NJ: Lawrence Erlbaum Associates.
- COP M. 1991. Collocations in the Bilingual Dictionary. *Pages 2775–2778 of: HAUSMANN F. J., REICHMANN O., WIEGAND H. E., & ZGUSTA L. (eds), Wörterbücher. Dictionaries. Dictionnaires. Ein internationales Handbuch zur Lexicographie. An international Encyclopedia of Lexicography. Encyclopédie internationale de lexicographie*, vol. III. Berlin - New York. Walter De Gruyter.
- COSERIU E. 1966. Structure lexicale et enseignement du vocabulaire. In: *Actes du premier colloque international de linguistique appliqué*.
- COSERIU E. 1967. Lexikalische Solidaritäten. *Poetica*, 293–303. trad. it. “Solidarietà lessicali” in Coseriu, E. *Teoria del linguaggio e linguistica generale. Sette studi*, Bari, Laterza, 303–316, 1971.
- CUTLER A. 1982. Idioms: the older the colder. *Linguistic Inquiry*, 13(2), 317–320.
- D’ADDIO W. 1974. La posizione dell’aggettivo italiano nel gruppo nominale. *Pages 79–103 of: Fenomeni morfologici e sintattici nell’italiano contemporaneo. Atti del sesto congresso internazionale di studi*. Roma: Bulzoni.
- DAILLE B. 1994. *Approche mixte pour l’extraction automatique de terminologie: statistiques lexicales et filtres linguistiques*. Ph.D. thesis, Université Paris 7.
- DAILLE B. 1996. Study and implementation of combined techniques for automatic extraction of terminology. *Chap. 3, pages 49–66 of: KLAVANS J. L., & RESNIK P. (eds), The Balancing Act*. Cambridge, MA: MIT Press.
- DAMERAU F. J. 1971. *Markov Models and Linguistic Theory*. The Hague: Mouton.
- DE MAURO T., & VOGHERA M. 1996. *Scala Mobile. Un punto di vista sui lessemi complessi*. Bulzoni Editore. In Benincà P., Cinque G., De Mauro T., Vincent N. (a cura di), *Italiano e dialetti nel tempo*.
- DE MAURO T. 1980. *Guida all’uso delle parole*. Roma: Editori Riuniti.

- DE MAURO T. 1999-2007. *GRADIT - Grande Dizionario Italiano dell'Uso*. UTET.
- DE MAURO T. 2005. *La fabbrica delle parole*. Torino: UTET.
- DE MAURO T., & CHIARI I. 2005. *Parole e numeri. Analisi quantitative dei fatti di lingua*. Roma, Aracne.
- DE MAURO T., MANCINI F., VEDOVELLI M., & VOGHERA M. 1993. *LIP - Lessico di frequenza dell'italiano parlato*. Milano, Etaslibri.
- DELL'ORLETTA F. 2009. Ensemble system for part-of-speech tagging. *In: Proceedings of Evalita '09, Evaluation of NLP and Speech Tools for Italian*.
- DI SCIULLO A. M., & WILLIAMS E. 1987. *On the definition of word*. MIT Press, Cambridge, Massachusetts.
- DIAS G., GUILLORÉ S., & LOPES J. G. P. 1999. Language independent automatic acquisition of rigid multiword units from unrestricted text corpora. *In: Proceedings of Traitement Automatique des Langues Naturelles*.
- DUMAIS S. T. 2005. Latent Semantic Analysis. *Annual Review of Information Science and Technology*, **38**(188).
- DUNNING T. E. 1993. Accurate methods for the statistics of surprise and coincidence. *Computational Linguistics*, **19**(1), 61–74.
- EVERT S. 2004. *The Statistics of Word Cooccurrences: Word Pairs and Collocations*. Ph.D. thesis, Institut für maschinelle Sprachverarbeitung, University of Stuttgart.
- EVERT S. 2008. Corpora and collocations. *In: LÜDELING A., & KYTÖ M. (eds), Corpus Linguistics. An International Handbook*. Mouton de Gruyter, Berlin.
- EVERT S., & KRENN B. 2001. Methods for the qualitative evaluation of lexical association measures. *Pages 188–195 of: Proceedings of the 39th Annual Meeting of the Association for Computational Linguistics*.
- FAZLY A., COOK P., & STEVENSON S. 2009. Unsupervised type and token identification of idiomatic expressions. *Computational Linguistics*, **35**(1), 61–103.
- FAZLY A., & STEVENSON S. 2007 (June). Distinguishing Subtypes of Multiword Expressions Using Linguistically-Motivated Statistical Measures. *Pages 9–16 of: Proceedings of the Workshop on A Broader Perspective on Multiword Expressions*.
- FILLMORE C. J., KAY P., & O'CONNOR M. C. 1988. Regularity and idiomaticity in grammatical constructions: the case of LET ALONE. *Language*, 501–538.
- FIRTH J. R. 1957. *Papers in Linguistics 1934 - 1951*. Oxford University Press.

- FRANCIS W. N., & KUCERA H. 1964. *Manual of Information to accompany A Standard Corpus of Present-Day Edited American English, for use with Digital Computers*. Tech. rept. Department of Linguistics, Brown University, Providence, Rhode Island. Revised and amplified 1979.
- FRASER B. 1970. Idioms within a transformational grammar. *Foundations of Language*, **6**(1), 22–42.
- FREGE G. 1884. *Die Grundlagen der Arithmetik. Eine logisch-mathematische Untersuchung über den Begriff der Zahl*. W. Koebner, Breslau.
- GIULIANO V. E. 1965. Postscript: A personal reaction to the reading of the conference manuscript. *Pages 259–260 of: STEVENS M. E., GIULIANO V. E., & HEILPRIN L. B. (eds), Proceedings of the Symposium on Statistical Association Methods for Mechanized Documentation*.
- GLEDHILL C. 2000. *Collocations in Science Writing*. Tübingen, Gunter Narr.
- GRAFFI G. 2010. *Due secoli di pensiero linguistico: dai primi dell'Ottocento a oggi*. Roma: Carocci.
- GRANDI N. 2006. Considerazioni sulla definizione e la classificazione dei composti. *Annali Online di Ferrara - Lettere*, **I**(1), 31–52. <http://annali.unife.it/lettere/2006vol1/grandi.pdf>.
- GRECIANO G. 1983. *Signification et denotation en allemand: la sémantique des expressions idiomatiques*. numero monografico di *Recherches Linguistiques*, IX.
- GREIMAS A. J. 1960. Idiotismes, proverbes, dictions. *Cahiers de Lexicologie*, 41–61.
- GROSS M. 1981. Les bases empiriques de la notion de prédicat sémantique. *Langages*, 7–52.
- GROSS M. 1984. Une classification des phrases “figées” du français. In: ALTAL P., & MULLER C. (eds), *De la syntaxe à la pragmatique*. Amsterdam, Benjamins.
- HALLIDAY M. A. K. 1961. Categories of the Theory of Grammar. *Word*, **17**(3).
- HALLIDAY M. A. K. 1966. Lexis as a linguistic level. *Journal of Linguistics*, **2**(1), 57–67.
- HARRIS Z. S. 1954. Distributional Structure. *Word*, **10**(23), 146–162.
- HARRIS Z. S. 1968. *Mathematical Structures of Language*. New York: Wiley.

- HEID U., & GOJUN A. 2012. Term candidate extraction for terminography and CAT: and overview of TTC. *In: Proceedings of the 15th Euralex International Congress.*
- HEID U., & WELLER M. 2010. Corpus-derived data on German multiword expressions for lexicography. *In: Proceedings of the Euralex International Congress 2010.* [CD-ROM].
- HOCKETT C. F. 1956. *Idiom formation.* The Hague, Mouton. in M. Halle et al., For Roman Jakobson.
- HOEY M. 2000. A world beyond collocation: new perspectives on vocabulary teaching. *In: LEWIS M. (ed), Teaching collocations.* Hove: Language Teaching Publications.
- HOWARTH P. 1996. *Phraseology in English academic writing: some implications for language learning and dictionary making.* Niemeyer, Tübingen.
- JACKENDOFF R. 1997. *The architecture of the language faculty.* MIT Press.
- JUSTESON J. S., & KATZ S. 1995. Technical terminology: Some linguistic properties and an algorithm for identification in text. *Natural Language Engineering*, 1, 9–27.
- KAALEP H.-J., & MUISCHNEK K. 2003. Inconsistent selectional criteria in semi-automatic multi-word unit extraction. *Pages 27–36 of: Proceedings of the 7th Conference on Computational Lexicography and Text Research (COMPLEX 2003).*
- KAGEURA K., & UMINO B. 1996. Methods of automatic term recognition. *Terminology*, 3(2), 259–289.
- KATZ J. J. 1966. *The philosophy of language.* Harper and Row, London.
- KATZ J. J., & FODOR J. A. 1963. The structure of a semantic theory. *Language*, 170–210.
- KATZ J. J., & POSTAL P. M. 1963. Semantic interpretation of idioms and sentences containing them. *Quarterly Progress Report of MIT Research Laboratory of Electronics*, 257–262.
- KILGARRIFF A., & GREFENSTETTE G. 2001. Web as a Corpus. *In: Proceedings of Corpus Linguistics 2001.*
- KILGARRIFF A., RYCHLY P., SMRZ P., & TUGWELL D. 2004. The Sketch Engine. *Pages 105–116 of: Proceedings of EURALEX 2004.*

- KRENN B. 2000. *The Usual Suspects: Data-Oriented Models for the Identification and Representation of Lexical Collocations. Volume 7 of Saarbrücken Dissertations in Computational Linguistics and Language Technology*. Saarbrücken, Germany: DFKI and Universität des Saarlandes.
- KRENN B., EVERT S., & ZINSMEISTER H. 2004. Determining intercoder agreement for a collocation identification task. *Pages 89–96 of: Proceedings of KONVENS 2004*.
- LAMBRECHT K. 1984. Formulaicity, Frame Semantics, and Pragmatics in German binomial expressions. *Language*, 753–796.
- LAPATA M., McDONALD S., & KELLER F. 1999. Determinants of adjective-noun plausibility. *Pages 30–36 of: Proceedings of the 9th Conference of the European Chapter of the Association for Computational Linguistics (EACL 1999)*.
- LEMNITZER L. 1998. Komplexe lexikalische Einheiten in Text und Lexicon. *Pages 85–92 of: HEYER G., & WOLFF C. (eds), Linguistik und neue Medien*. Wiesbaden: DUV.
- LENCI A. 2009. Spazi di parole: metafore e rappresentazioni semantiche. *Paradigmi*, 83–100.
- LEPSCHY G. C. 1989. Lessico. *In: Sulla linguistica moderna*. Bologna, Il Mulino.
- LEXICAL COMPUTING LTD. 2014. *Statistics used in the Sketch Engine*. <http://trac.sketchengine.co.uk/raw-attachment/wiki/SkE/DocsIndex/ske-stat.pdf>.
- LIN D. 1998. Extracting collocations from text corpora. *Pages 57–63 of: Proceedings of the First Workshop on Computational Terminology*.
- LIN D. 1999. Automatic identification of non-compositional phrases. *In: Proceedings of ACL99*.
- LO CASCIO V. 2011. *Dizionario Combinatorio Italiano*. John Benjamins.
- LYDING V., STEMLE E., BORGHETTI C., BRUNELLO M., CASTAGNOLI S., DELL'ORLETTA F., DITTMANN H., LENCİ A., & PIRRELLI V. 2014. The PAISÀ Corpus of Italian Web Texts. *Pages 36–43 of: Proceedings of the 9th Web as Corpus Workshop (WaC-9)*. Gothenburg, Sweden: Association for Computational Linguistics.
- LYONS J. 1977. *Semantics, 2 Vol.* Cambridge, Cambridge University Press.
- MAKKAI A. 1972. *Idiom structure in English*. The Hague, Mouton.

- MALKIEL Y. 1959. Studies in irreversible binomials. *Lingua*, 113–160.
- MANNING C., & SCHÜTZE H. 1999. *Foundations of Statistical Natural Language Processing*. Cambridge, MIT Press.
- MARTINET A. 1967. Syntagme et syntème. *La Linguistique*, 1–14.
- MASINI F. 2007. *Parole sintagmatiche in italiano*. Ph.D. thesis, Università degli Studi di Roma Tre.
- MASINI F. 2008. Binomi coordinati in italiano. *Pages 563–571 of: CRESTI E. (ed), Prospettive nello studio del lessico italiano. Atti SILFI 2006*, vol. II. Firenze: Università di Firenze.
- MASINI F. 2009. Combinazioni di parole e parole sintagmatiche. *Pages 191–209 of: LOMBARDI VALLAURI E., & MEREU L. (eds), Spazi linguistici. Studi in onore di Raffaele Simone*. Roma: Bulzoni.
- MASINI F., & THORNTON A. M. 2007 (September). Italian VeV lexical constructions. *In: BOOIJ G., RALLI A., & SCALISE S. (eds), Proceedings of the 6th Mediterranean Morphology Meeting*.
- MCCARTHY D., KELLER B., & CARROLL J. 2003. Detecting a Continuum of Compositionality in Phrasal Verbs. *Pages 73–80 of: Proceedings of ACL-SIGLEX Workshop on Multiword Expressions*.
- MCINTOSH C., FRANCIS B., & POOLE R. 2009. *Oxford Collocations Dictionary for students of English*. Oxford University Press.
- MEL'ČUK I. A. 1998. *Collocations and Lexical Functions*. Oxford: Clarendon Press. In Cowie A. P. (a cura di), *Phraseology, Theory, Analysis and Application*.
- MIGLIORINI B. 1960. *Storia della lingua italiana*. Milano: Bompiani.
- MILLER G. A., & CHARLES W. G. 1991. Contextual Correlates of Semantic Similarity. *Language and Cognitive Processes*, **VI**, 1–28.
- MONTI J., BARREIRO A., ELIA A., MARANO F., & NAPOLI A. 2011. Taking on new challenges in multi-word unit processing for machine translation. *Pages 11–19 of: Second International Workshop on Free/Open-Source Rule-Based Machine Translation. Barcelona, 20-21 January 2011*.
- MOON R. 1997. Vocabulary Connections: Multi-Word Items in English. *In: SCHMITT N., & MCCARTHY M. (eds), Vocabulary: Description, Acquisition and Pedagogy*. Cambridge, Cambridge University Press.
- NELSON M. 2000. *A corpus-based study of the lexis of business English and business English teaching materials*. Ph.D. thesis, University of Manchester.

- NERIMA L., SERETAN V., & WEHRLI E. 2003. Creating a multilingual collocation dictionary from large text corpora. *Pages 131–141 of: Companion Volume to the Proceedings of the 10th Conference of the European Chapter of the Association for Computational Linguistics.*
- NEWMAYER F. J. 1972. The insertion of idioms. *In: Papers from the 8th regional meeting of the Chicago Linguistic Society.* Chicago, Chicago Linguistic Society.
- NEWMAYER F. J. 1974. The regularity of idiom behavior. *Lingua*, 327–342.
- NISSIM M., & ZANINELLO A. 2011. A quantitative study on the morphology of Italian multiword expressions. *Lingue e Linguaggio*, 283–300.
- NIVRE J. 2005. *Dependency grammar and dependency parsing.* Tech. rept. Växjö University.
- NORVIG P. 2012. Colorless green ideas learn furiously: Chomsky and the two cultures of statistical learning. *Significance. Special Issue: Big Data*, 9(4), 30–33.
- NUNBERG G. 1978. *The pragmatics of reference.* Bloomington, Indiana University Linguistics Club.
- PAUL H. 1880. *Prinzipien der Sprachgeschichte.* Halle, Niemeyer.
- PEREIRA F. 2002. Formal grammar and information theory: together again? *In: NEVIN B. E., & JOHNSON S. M. (eds), The Legacy of Zellig Harris.* Benjamins.
- PEREIRA L., STRAFELLA E., DUH K., & MATSUMOTO Y. 2014. Identifying collocations using cross-lingual association measures. *Pages 109–113 of: Proceedings of the 10th Workshop on Multiword Expressions (MWE 2014).*
- PIKE K. L. 1967. *Language in relation to a unified theory of the structure of human behavior.* The Hague, Mouton.
- PORZIG W. 1934. Wesenhafte Bedeutungsbeziehungen. *Beträge zur Geschichte der Deutschen Sprache und Literatur*, 70–97.
- PRENDERGAST T. 1864. *The mastery of languages; or the art of speaking foreign tongues idiomatically.* Richard Bentley, London.
- PUSTEJOVSKY J. 1995. *The Generative Lexicon.* Cambridge, The MIT Press.
- RAMAT P. 1990. Definizione di “parola” e sua tipologia. *In: BERRETTA M., MOLINELLI P., & VALENTINI A. (eds), Parallela 4. Morfologia.* Gunter Narr, Tübingen.
- RAMAT P. 2005. Per una definizione di parola. *In: Pagine linguistiche.* Laterza.

- RAMISCH C., VILLAVICENCIO A., MOURA L., & IDIART M. 2008. Picking them up and Figuring them out: Verb-Particle Constructions, Noise and Idiomaticity. *Pages 49–56 of: Proceedings of the 12th Conference on Computational Natural Language Learning (CoNLL 2008).*
- RAMISCH C., MEDEIROS CASELI H., VILLAVICENCIO A., MACHADO A., & FINATTO M. J. 2010a. A Hybrid Approach for Multiword Expression Identification. *In: Proceedings of the 9th International Conference on Computational Processing of Portuguese Language (PROPOR 2010).*
- RAMISCH C., VILLAVICENCIO A., & BOITET C. 2010b. mwetoolkit: a Framework for Multiword Expression Identification. *In: Proceedings of the LREC Workshop on Exploitation of multilingual resources and tools for Central and (South) Eastern European Languages.*
- RUNDELL M. 2002. *Macmillan English Dictionary for Advanced Learners.* Macmillan.
- RUWET N. 1983. Du bon usage des espressions idiomatiques dans l'argumentation en syntaxe générative. *Revue Québécoise de linguistique*, **13**(1), 9–145.
- SAG I. A., BALDWIN T., BOND F., COPESTAKE A., & FLICKINGER D. 2001. Multiword Expressions: A Pain in the Neck for NLP. *Pages 1–15 of: Proceedings of the 3rd International Conference on Intelligent Text Processing and Computational Linguistics (CICLing-2002).*
- SAUSSURE F. D. 1922. *Cours de linguistique générale. A cura di C. Bally e A. Sechehaye.* Payot, Paris.
- SCALISE S. 1994. *Le strutture del linguaggio.* Bologna, Il Mulino.
- SCALISE S., & BISETTO A. 2008. *La struttura delle parole.* Bologna, Il Mulino.
- SCHENK A. 1995. The syntactic behavior of idioms. *In: EVERAERT M., VAN DER LINDEN E.-J., SCHENK A., & SCHREUDER R. (eds), Idioms: structural and psychological perspectives.* Hillsdale, Lawrence Erlbaum.
- SCHMID H. 1994. Probabilistic Part-of-Speech Tagging Using Decision Trees. *In: Proceedings of International Conference on New Methods in Language Processing, Manchester UK.*
- SCHONE P., & JURAFSKY D. 2001. Is Knowledge-free Induction of Multiword Unit Dictionary Headwords a Solved Problem? *Pages 100–108 of: Proceedings of the 2001 Conference on Empirical Methods in Natural Language Processing.*
- SCHÜTZE H. 1998. Automatic Word Sense Discrimination. *Computational Linguistics*, **24**(1), 97–123.

- SEARLE J. R. 1975. *Indirect speech acts*. in Cole, P. e Morgan, J. L. (eds.). *Speech Acts*. Academic Press, New York.
- SECHEHAYE A. 1921. Locutions et composés. *Journal de Psychologie Normale et Pathologique*, 654–675.
- SERETAN V. 2011. *Syntax-based Collocation Extraction*. Berlin, Germany: Springer.
- SERIANNI L. 1989. *Grammatica italiana*. Torino: UTET.
- SHANNON C. 1948. A Mathematical Theory of Communication. *The Bell System Technical Journal*, **27**, 379–423.
- SHANNON C., & WEAVER W. 1949. *A Mathematical Model of Communication*. Urbana, IL, University of Illinois Press.
- SIMONE R. 1990. *Fondamenti di linguistica*. Roma - Bari, Laterza.
- SINCLAIR J. 1966. Beginning the study of lexis. In: BAZELL C. E., CATFORD J., HALLIDAY M. A. K., & ROBINS R. H. (eds), *In Memory of J. R. Firth*. Longman.
- SINCLAIR J. 1993. Text Corpora: Lexicographer's needs. *Zeitschrift für Anglistik und Americanistik*, **XLI**(1), 5–13.
- SINCLAIR J., & RENOUF A. 1991. Collocational framework in English. In: AIJMER K., & ALTENBERG B. (eds), *English Corpus Linguistics. Studies in honour of Jan Svartvik*. London and New York, Longman.
- SINCLAIR J. 1991. *Corpus, concordance, collocation*. Oxford University Press.
- SMADJA F. 1993. Retrieving collocations from text: Xtract. *Computational Linguistic*, **19**(1), 143–177.
- SMADJA F., MCKEOWN K. R., & HATZIVASSILOGLU V. 1996. Translating collocations for bilingual lexicon: a statistical approach. *Computational Linguistics*, **22**(1), 1–38.
- SOMERS H. 2003. *Computers and Translation: a Translator's Guide*. John Benjamins, Amsterdam.
- SQUILLANTE L. 2014. Towards and Empirical Subcategorization of Multiword Expressions. *Pages 77–81 of: Proceedings of the 10th Workshop on Multiword Expressions (MWE 2014)*.

- STEVENS M. E., GIULIANO V. E., & HEILPRIN L. B. (eds). 1965. *Proceedings of the Symposium on Statistical Association Methods for Mechanized Documentation*. Vol. 269.
- STUBBS M. 1995. Collocations and semantic profiles: On the cause of the trouble with quantitative studies. *Functions of Language*, 1, 23–55.
- SWEET H. 1891. *A new English grammar, logical and historical*. Oxford Clarendon Press.
- TAPANAINEN P., PIITULAINEN J., & JÄRVINEN T. 1998. Idiomatic object usage and support verbs. In: *36th Annual Meeting of the Association for Computational Linguistics*.
- TENSIÈRE L. 1959. *Éléments de syntaxe structurale*. Editions Klincksieck.
- TIBERII P. 2012. *Dizionario delle Collocazioni. Le combinazioni delle parole in italiano*. Zanichelli.
- URZÌ F. 2009. *Dizionario delle Combinazioni Lessicali*. Convivium, Lussemburgo.
- VAN DE CRUYS T., & VILLADA MOIRÓN B. 2007. Semantics-based Multiword Expression Extraction. In: *Proceedings of the ACL 2007 Workshop on Multiword Expressions: A Broader Perspective*.
- VAN DER WOUDE T. 1997. *Negative contexts. Collocations, polarity and multiple negation*. Routledge.
- VAN ROEY J. 1990. *French-English Contrastive Lexicology. An Introduction*. Louvain-La-Neuve. Peeters.
- VENKATAPATHY S., & JOSHI A. K. 2005. Measuring the relative compositionality of verb-noun (V-N) collocations by integrating features. *Pages 899–906 of: Proceedings of HLT-EMNLP 05*.
- VIETRI S. 1985. *Lessico e sintassi delle espressioni idiomatiche. Una tipologia tassonomica dell'italiano*. Liguori Editore.
- VILLADA MOIRÓN B., & TIEDEMANN J. 2006. Identifying idiomatic expressions using automatic word-alignment. *Pages 33–40 of: Proceedings of the EACL 2006 Workshop on Multiword Expressions in a Multilingual Context*.
- VINCENT N. 1986. La posizione dell'aggettivo in italiano. *Pages 181–195 of: Tema-rema in italiano. Symposium*. Frankfurt Am Main: Gunter Narr Verlag.
- VOGHERA M. 1994. Lessemi complessi: percorsi di lessicalizzazione a confronto. *Lingua e stile*, XXIX(2).

VOGHERA M. 2004. Polirematiche. *Pages 56–69 of: GROSSMANN M., & RAINER F. (eds), La formazione delle parole in italiano.* Tübingen: Niemeyer.

WASOW T., SANG I., & NUNBERG G. 1983. *Idioms: an interim report.* in Hattori S., Inoue K. (cur.), *Proceedings of the 13th international congress of linguists.* The Hague, Mouton.

WEINRICH U. 1966. *Exploration in semantic theory.* in Sebeok, T. S., *Current trends in linguistics III: theoretical foundations.* The Hague, Mouton.

WEINRICH U. 1969. *Problems in the analysis of idioms.* in Puhvel, J., *Substance and structure of language.* Berkely, University of California Press.

WELLER M., & FRITZINGER F. 2010. A Hybrid Approach for the Identification of Multiword Expressions. *In: Proceedings of the SLTC 2010 Workshop on Compounds and Multiword Expressions.*

WERMTER J., & HAHN U. 2004. Collocation Extraction Based on Modifiability Statistics. *Pages 980–986 of: Proceedings of COLING 2004.*